

2ej  
21



Universidad Nacional Autónoma de México

FACULTAD DE CIENCIAS

SELECCION DE VARIABLES EN ANALISIS  
DE REGRESION BAYESIANO.

T E S I S  
Que para obtener el título de  
M A T E M A T I C O  
P R E S E N T A  
HORTENSIA REYES CERVANTES

México, D. F.

1987



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.

# I N D I C E

	PÁG.
INTRODUCCIÓN.	1
CAPITULO I. ANALISIS DE REGRESION .	1
CAPITULO II. METODOS DE SELECCION DE VARIABLES.	22
CAPITULO III. CONCLUSIONES	60
ANEXO.	70
BIBLIOGRAFIA.	73

## INTRODUCCIÓN

El origen de esta tesis surge del interés de conocer el material referente al Análisis de Regresión Bayesiano, orientado a resolver el problema de selección de variables, "Ya que el Análisis de Regresión es una de las técnicas estadísticas más empleadas para analizar múltiples factores de datos" (Montgomery, 1982, p.v)

Es importante tener en mente, que no siempre es deseable efectuar una selección de variables ya que de realizarse pueden cometerse errores que llevan a perder información y/o en algunas ocasiones se complique innecesariamente el problema y esto, al final de cuentas, repercutirá en lo bueno que sean los resultados finales. Por ejemplo, tomando como referencia el procedimiento de selección dado por Lindley, (1968), se tendrían que tomar en cuenta: costos, captación y riesgos de haber elegido a las variables explicativas. Además, el peligro de una mala inferencia se implementa por la selección de variables.

El no haber hecho una selección de variables, implica el resolver el problema de regresión con todas las variables explicativas que se tengan y esto, actualmente ya no resulta un trabajo difícil en algunos casos, pues en gran medida han ayudado los avances tecnológicos de las nuevas computadoras.

En general, la elección del mejor modelo de regresión no es un mecanismo estadístico único (Weisberg, 1980, p.174), la razón se debe a la diversidad de las situaciones en que se desarrollan los problemas de regresión; y a la manera distinta que cada individuo tiene de concebir y actuar ante un problema.

La elección del mejor modelo de regresión debe establecer un compromiso entre dos criterios (opuestos) que son presentados por Draper-Smith (1981, p.294), las cuales son:

1. Incluir en el modelo, todas las posibles variables explicativas para predecir mejor a la variable dependiente.
2. Dado que cada variable explicativa requiere un costo de obtención y captación de información, se pretende elegir un menor número de variables explicativas.

La justificación para llevar a cabo una selección de variables es ocasionada por alguna o todas las siguientes causas, que se pueden presentar en un problema de regresión.

- 1º. El problema de regresión cuenta con "Muchas Variables", por lo que debido a las consideraciones hechas por el interesado y/o estadístico, se llega a la conclusión que hay información repetitiva en las observaciones.
- 2º. Las características propias del problema de regresión hacen que el número de variables explicativas sea menor que el número de observaciones.
- 3º. Debido al costo desmedido en el manejo y captación de las variables explicativas, se desea que para metas futuras de mantenimiento en el modelo de regresión, éste cuente con el menor número de variables explicativas.

A continuación se presenta de manera breve el contenido del trabajo.

En el Capítulo 1, se plantea el Modelo de Regresión Bayesiana se dan algunas indicaciones en donde el Modelo de Regresión Clásico y el Bayesiano coinciden. En el Capítulo 2, se presentan algunos de los procedimientos de Selección Bayesiano y en el Capítulo 3 se presentan las conclusiones, así como comentarios respecto a cual podría ser la mejor técnica de selección.

**ANALISIS  
DE  
REGRESION**

## CAPITULO I

### ANALISIS DE REGRESION.

El Análisis de Regresión es una de las técnicas estadísticas más empleadas para analizar situaciones en donde el objetivo esencial es establecer una función que permita al estadístico relacionar un conjunto de variables llamadas independientes  $(X_1, \dots, X_n)$  con otra variable llamada dependiente  $(Y)$ .

Los objetivos principales de los modelos de regresión son a grandes rasgos: 1) Descripción de la relación funcional entre la variables independiente y la variable dependiente, y 2) Predicción de la variable dependiente.

Existen dos enfoques estadísticos alternativos que solventan los problemas de regresión, estos son: el punto de vista bayesiano y el punto de vista clásico. El enfoque que se utilizará a lo largo de este trabajo, es el bayesiano, luego entonces el enfoque clásico (Draper-Smith, 1981) no se desarrollará y únicamente se referirá a él, en casos especiales para comparar algunos resultados.

En el Análisis de Regresión Bayesiano, el estadístico puede introducir el conocimiento disponible sobre el problema, a la información producida por las observaciones, mediante el Teorema de Bayes, y además, se puede proporcionar un criterio de decisión basándose en una función de pérdida o ganancia. Si desde un enfoque bayesiano se utiliza una distribución inicial de referencia, entonces las similitudes entre los resultados obtenidos por ambos enfoques son gran



des.

En lo referente a problemas de inferencia, pueden presentarse situaciones en que no haya necesidad de elegir formalmente un estimador porque no se tiene o no se quiere utilizar una función de pérdida, ya que lo único que se pretende es hacer la Inferencia Estadística con respecto a los parámetros involucrados. Por ejemplo, expresar el conocimiento disponible sobre el parámetro de localización y/o el parámetro de escala, así como otras inferencias acerca de las distribuciones de las observaciones.

En este capítulo, se pretende dar un panorama general del Análisis de Regresión Bayesiano desde dos puntos de vista o necesidades del problema: Inferencia y Decisión. Por último se tocará muy sucintamente el tema del Modelo Lineal General.

## 1. MODELO DE REGRESIÓN LINEAL.

La formulación del problema de Regresión Bayesiano es necesariamente más complicada y requiere de más cuidado que la descripción técnica muestral clásica. Esta diferencia se debe a 1) que todas las cantidades  $\{y, x_1, \dots, x_n\}$  y los parámetros son variables aleatorias, y 2) todas las afirmaciones probabilísticas son condicionales a los eventos de interés (Lindley, 1968).

El problema de Regresión Lineal Múltiple está contemplado en el área que involucra los modelos lineales, porque se supone que la media de cada una de las observaciones es una combinación lineal de un conjunto de parámetros fijos. Por lo cual, se puede escribir de la siguiente forma

$$y = X\theta + \epsilon \quad (1.1)$$

donde:  $\theta$  es un vector de  $(k+1)$  parámetros desconocidos cuyos valores se encuentran en  $\mathcal{R}$ ,  $\epsilon$  es un vector de  $n \times 1$  con distribución Normal Multivariada con media cero y varianza común  $\sigma^2$  e independientes  $\{\epsilon_i \sim N(0, \sigma^2) \forall i=1, \dots, n\}$ , por lo tanto  $y$  es un vector  $(n \times 1)$  que tiene distribución Normal para toda  $i=1, \dots, n$  ( $y = [y_1, \dots, y_n]^T$ ); cuya media es una combinación de  $k$  componentes de  $\theta$  ( $E(y/\theta) = X\theta$ ) y cuya varianza es  $\sigma^2$  ( $V(y/X) = \sigma^2 I$ ). Por último,  $X$  matriz  $(n \times k)$  de valores  $x_{ij}$  (con  $i=1, \dots, n$ ;  $j=1, \dots, k$ ) cuyo rango es igual a  $k$ ; que pueden representar niveles o funciones de una variable de un experimento, o quizá se representen valores de variables indicadoras (representados por unos o ceros).

Se denota por  $\sigma^2$  a la varianza cuyo valor debe ser positivo, y

se puede demostrar que sus estadísticas respectivas  $\theta$  y  $s^2$  son independientes. Estas estadísticas más adelante se presentarán.

La información expresada en la muestra, es decir, las observaciones  $X_1, \dots, X_n$  se representan mediante la función de  $\underline{y}$ , que es crita en forma matricial es

$$p(\underline{y}/\theta, \sigma^2, X) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \sigma^{-n} e^{-\frac{1}{2\sigma^2}(\underline{y}-X\theta)^t(\underline{y}-X\theta)} \quad (1.2)$$

Haciendo una breve observación sobre esta densidad, se tiene que tomando  $w = \frac{1}{\sigma^2}$ , la precisión, se sigue conservando la normalidad en la densidad de  $\underline{y}$ . Solo se obtiene una expresión que para algunas personas proporciona un manejo más cómodo en cuanto a su interpretación.

$$p(\underline{y}/\theta, w, X) = \left(\frac{1}{\sqrt{2\pi}}\right)^n w^{n/2} e^{-\frac{w}{2}(\underline{y}-X\theta)^t(\underline{y}-X\theta)} \quad (1.2')$$

Retomando la expresión (1.2), esta se puede expresar como

$$p(\underline{y}/\theta, \sigma^2, X) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \sigma^{-n} e^{-\frac{1}{2\sigma^2}[(\underline{y}-\hat{\underline{y}})^t(\underline{y}-\hat{\underline{y}}) + (\theta-\hat{\theta})^t X^t(\theta-\hat{\theta})]} \quad (1.3)$$

que consiste en sumar y restar las estadísticas  $\hat{\theta}$ ,  $s^2$ , agrupando adecuadamente. Este procedimiento es el mismo que Zellner hizo (1968, p.66) y así obtuvo la densidad de  $\underline{y}$  como

$$p(\underline{y}/\theta, \sigma^2, X) = \left(\frac{1}{\sqrt{2\pi}}\right)^n \sigma^{-n} e^{-\frac{1}{2\sigma^2}[\nu s^2 + (\theta-\hat{\theta})^t X^t(\theta-\hat{\theta})]} \quad (1.4)$$

con

$$\hat{\theta} = (X^T X)^{-1} X^T y, \quad (1.5)$$

$$s^2 = \frac{1}{v} (y - \hat{y})^T (y - \hat{y}), \quad \hat{y} = X \hat{\theta}, \quad v = n - k \quad (1.6)$$

A partir de estas expresiones y apoyados en otros resultados se puede demostrar lo siguiente: 1)  $(\hat{\theta}, s^2)$  es una estadística suficiente conjunta de  $(\theta, \sigma^2)$ , 2) que el estimador  $(\hat{\theta}, \hat{\sigma}^2)$ , tiene una distribución Normal Multivariada con varianza  $\sigma^2 (X^T X)^{-1}$  (esto es,  $\hat{\theta} \sim N_k(\theta, \sigma^2 (X^T X)^{-1})$ ) y 3)  $v s^2$  tiene una distribución independiente de  $\hat{\theta}$  con distribución ji-cuadrada con  $v$ -grados de libertad (se denota por  $\frac{v s^2}{\sigma^2} \sim \chi_v^2$ ).

Después de haber perdido normalidad en las observaciones y haber obtenido su densidad,  $p(y|\theta, \sigma^2)$  es necesario establecer en seguida que tipo de distribución inicial (referente al parámetros desconocido) es la contemplada por el interesado y/o estadístico. Para lo cual se detallan a continuación los casos de 1.1) cuando el interesado expresa su conocimiento acerca del (os) parámetro(s) desconocido(s) y aplica la regla de Bayes. Concluyendo con una distribución posterior acerca de(los) parámetro(s),  $\theta$  cuando 1.2) existen situaciones en que el interesado no expresa su conocimiento por no poder o querer, ocasionando así que no se establezca una distribución inicial informativa acerca de los parámetros desconocidos.

### 1.1. DISTRIBUCIÓN INICIAL INFORMATIVA,

Esgrimiendo el hecho de que la información inicial se necesita

incorporar mediante una distribución adecuada, DeGroot (1970, p.183) propone utilizar el Modelo Normal Gamma, que trae como consecuencia el establecimiento de que la distribución condicional de  $\theta$  cuando  $W=w$  es una distribución normal multivariada con vector de medias  $\mu$  y matriz de precisión  $w\tau$  (esto es,  $\theta|W=w \sim N_k(\mu, w\tau)$ ), y que la distribución marginal de  $W$  es una Gamma con parámetros  $\alpha$  y  $\beta$  (ambos positivos). Entonces la distribución inicial conjunta de  $\theta$  y  $w$  es:

$$f(\theta, w) \propto w^{\alpha/2} e^{-\frac{w}{2}(\theta-\mu)'\tau(\theta-\mu)} w^{\alpha-1} e^{-\beta w} \quad (1.7)$$

Por lo tanto, la distribución posterior  $p(\theta, w|y, x)$  se obtiene haciendo el producto de (1.2') y (1.7), concluyendo con

$$p(\theta, w|y, x) \propto w^{\alpha/2} e^{-\frac{w}{2}(\theta-\mu)'\tau(\theta-\mu)} w^{\alpha-1} e^{-\beta w} e^{-\frac{w}{2}(y-x\theta)'\tau(y-x\theta)} \quad (1.8)$$

tomando a

$$\mu_1 = (\tau + x'x)^{-1} (\tau\mu + x'y) \quad y \quad (1.9)$$

$$\beta_1 = \beta + \frac{1}{2} [(y-x\mu_1)'\tau(y-x\mu_1) + (\mu-\mu_1)'\tau\mu] \quad \text{en (1.8)}$$

se obtiene la distribución posterior  $p(\theta, w|y, x)$ , así:

$$p(\theta, w|y, x) \propto w^{\alpha/2 + \frac{\alpha}{2} - 1} e^{-\frac{w}{2}(\theta-\mu_1)'\tau(x)(\theta-\mu_1) - \beta_1 w} \quad (1.10)$$

Las marginales de esta distribución involucradas son: La distribución condicional de  $\theta$  dado  $W=w$  que es una distribución normal multivariada con media  $\mu_1$  y matriz de precisión  $w(\underline{\Sigma} + \underline{X}^t \underline{X})$  (es decir,  $\theta|W=w \sim N_n(\mu_1, w(\underline{\Sigma} + \underline{X}^t \underline{X}))$ ), la distribución marginal de  $w$  es una Gamma con parámetros  $\alpha + \frac{n}{2}$  y  $\beta_1 (w \sim \text{Gamma}(\alpha + \frac{n}{2}, \beta_1))$  y la distribución marginal de  $\theta$  tiene una distribución  $t$ -multivariada con  $2\alpha+n$  grados de libertad, vector de localización  $\mu_1$  y matriz de precisión  $\frac{2\alpha+n}{2\beta_1} (\underline{\Sigma} + \underline{X}^t \underline{X})$ .

### 1.1.1. INFERENCIA.

Las inferencias acerca de los parámetros desconocidos son el resultado del análisis y estudio de la distribución posterior conjunta o marginal. Para realizar estos objetivos es de gran utilidad el obtener regiones de alta densidad.

Las regiones de alta densidad de tamaños  $(1-\alpha)$  son comunes en la Estadística Bayesiana (Box-Tiao, 1971, pp.121-136) y se construyen sobre el espacio paramétrico de interés. La caracterización de estas regiones  $(R)$ , está dada por cualquiera de las siguientes condiciones: 1) que la densidad de probabilidad de todo punto interior a  $R$  sea mayor a la densidad correspondiente a todo punto exterior (esto es, si  $\theta \in R, \theta_2 \notin R \Rightarrow p(\theta_1 | y) \geq p(\theta_2 | y)$ ); 2) la región  $R$ , para una probabilidad contenida  $1-\alpha$ , debe ser tal que tenga el volumen más pequeño entre todas las regiones en el espacio paramétrico con contenido de probabilidad  $1-\alpha$ , esta propiedad se expresa así: si  $A$  y  $R$  son conjuntos tales que  $P_r(\theta \in A | y) = P_r(\theta \in R | y) = 1-\alpha$ , entonces el volumen  $(A) \geq$  volumen  $(R)$ .

Además, se puede demostrar que al tomar una condición trae como consecuencia la otra.

Si definimos una transformación  $\Phi = f(\theta)$ , de los parámetros de  $\theta$  a  $\Phi$ , para alguna región contenida en el espacio de  $\theta$  y si la región de  $\theta$  es de alta densidad, su región transformada en  $\Phi$  no será de alta densidad, a menos que la transformación utilizada sea lineal.

Para determinar si un parámetro puntual  $\theta_0$ , se encuentra en una región,  $R_\alpha$ , de alta densidad que no está determinada. Se puede proceder de dos formas: 1) determinar analíticamente la región  $\theta$  2) calcular  $P_r \{p(\theta/y) > p(\theta_0/y)\} = \rho$  si  $\theta_0 \in R_\alpha$ , en caso contrario  $\theta_0 \notin R_\alpha$ .

En esta expresión, la función de densidad  $p(\theta/y)$  es tratada como una variable aleatoria, pues  $p(\theta_0/y)$  es constante.

### 1.1.2 DECISION.

Cuando se quiere estimar un parámetro desconocido, digamos  $\omega = (\omega_1, \dots, \omega_k)^t$  que por ejemplo puede ser el vector de medias o de precisión, a través de la teoría de decisión (Bernardo, 1980, pp 4.2a 4.26) es necesario especificar (de acuerdo a los axiomas de coherencia) los siguientes conjuntos: El conjunto de decisiones y que en este caso coincide con el espacio paramétrico  $\{\omega \in W\}$ , el conjunto de consecuencias (C) y construir una función de utilidad que describa las preferencias del decisor, entre las posibles consecuencias  $(U(d, \omega))$ .

Partiendo de la definición de utilidad de DeGroot (1971, pp.91-92) pueden definirse fácilmente las funciones de ganancia o de pérdida, que son resultados de una transformación apropiada de la función de utilidad. Esto es, si  $U$  es una función de utilidad, la transformación  $aU + b$  define una función de ganancia si  $a > 0$ ; y en el caso de que  $a < 0$  se obtiene una función de pérdida.

Retomando lo anterior, se ha definido el problema de decisión sobre el espacio paramétrico  $\omega$ , en donde la finalidad es seleccionar una decisión  $d^*$  (con  $d^* \in D$ ) que minimice o maximice la función de utilidad.

Si se cuenta con los axiomas de coherencia (Bernardo, 1980, pp.4.2 4.26) el único criterio para elegir alternativas consiste en maximizar la utilidad esperada sobre el espacio de decisiones. <sup>(1)</sup> Equivalen-

(1) Tomando  $p \in \mathcal{P}$  (la clase de distribuciones de probabilidad sobre el espacio paramétrico  $\Theta$ ), y que  $g$  es una función de utilidad de  $d$ , el criterio consiste en

$$\max_{d \in D} E(g|P) = \int_{\Theta} g(r) dP(r), \quad \text{si } E(g|P) \text{ existe}$$



temente se puede expresar el criterio de decisión como maximizar la ganancia esperada o minimizar la pérdida esperada.

Los ejemplos más comunes de funciones de pérdidas son:

La función de pérdida lineal

$$L(d, w) = a|w-d|, \quad a > 0 \quad (1.14)$$

y la cuadrática

$$L(d, w) = b(w-d)^2, \quad b > 0. \quad (1.15)$$

Tomando al parámetro  $w \in \mathbb{R}^k$  y  $d \in \mathbb{R}^k$ , la función de pérdida se define como  $L(d, w) = (w-d)^t A (w-d)$  donde  $A$  es una matriz simétrica. El objetivo en este caso es maximizar  $E[(w-d)^t A (w-d)]$  en  $d$ .

Algunas veces, como se verá más adelante, se puede presentar el problema de regresión como un planteamiento de Teoría de Decisiones.

## 1.2 DISTRIBUCIÓN INICIAL DE REFERENCIA.

Como ya se habla comentado si no se cuenta con una distribución inicial no es posible obtener la distribución posterior del parámetro desconocido. De aquí, que es necesario un mecanismo para sustituirlo matemáticamente en el Teorema de Bayes y a estas funciones se les llama Distribuciones Iniciales de Referencia.

Existen diversos métodos para encontrar Iniciales de Referencia, de las cuales sólo se referiran a las concluidas por Degroot (1970, pp.155-222), Box y Tiao (1971, pp. 1-74) y Bernardo (1980). Estos métodos en general no coinciden, sin embargo suponiendo normalidad en las observaciones y siendo de interés solo uno de los parámetros desconocidos ( $\theta$  ó  $w = 1/\sigma^2$ ) se tienen Distribuciones Iniciales de Referencia iguales.

Esta situación se plantea de la siguiente manera <sup>(1)</sup>: si el vector de medias ( $\theta$ ) es desconocido y la precisión es conocida, la inicial de referencia que se propone es:

$$p(\theta, w = 1/\sigma^2) \propto \text{constante.} \quad (1.16)$$

Ahora si el vector de medias ( $\theta$ ) y la precisión ( $w$ ) son desconocidos, y únicamente es de interés un parámetro, Bernardo (1979) propone la inicial de referencia dada por

$$p(w) \propto \frac{1}{w} \text{ y } p(\theta/w) \propto \text{constante} \quad (1.17)$$

(1) Los resultados coinciden con las reglas de Jeffreys (Zellner, pp.42-44), que son distribuciones iniciales que representan poco conocimiento o ignorancia.

A continuación se distinguen dos situaciones prácticas de incertidumbre acerca de los parámetros involucrados en el desarrollo del Análisis de Regresión: 1.2.1 varianza conocida y media desconocida y 1.2.2 varianza y media desconocidas.

### 1.2.1 $\theta$ DESCONOCIDA Y $\sigma^2$ CONOCIDA.

Caso poco habitual en la práctica, pero sirve de sustento teórico para el caso precedente (Sección 1.2.2).

Usando la regla de Bayes y expresando el resultado en términos de las estadísticas suficientes, para obtener la distribución final del parámetro desconocido, se tiene

$$p(\theta/y, x) = p(\hat{\theta}/y, x, \theta) p(\theta/\hat{\theta}, y, x)$$

$$p(\theta/y, x) \propto p(\hat{\theta}/\theta) p(\theta), \quad (1.18)$$

se substituyen los términos (1.4) y (1.16) en (1.18) y quitando lo que no es desconocido, se tiene

$$p(\theta/y, x) \propto e^{-\frac{1}{2\sigma^2}[(\theta - \hat{\theta})^2 x^2 x(\theta - \hat{\theta})]}, \quad (1.19)$$

encontrando la constante de integración sobre  $\theta$

$$\int_{\theta_1} \dots \int_{\theta_n} e^{-\frac{1}{2\sigma^2}[(\theta - \hat{\theta})^2 x^2 x(\theta - \hat{\theta})]} d\theta_1 \dots d\theta_n = \frac{\sigma^n (\sqrt{2\pi})^n}{(x'x)^{1/2}}, \quad (1.20)$$

por lo que substituyendo (1.20) en (1.19), se encuentra

$$p(\underline{\theta} | \underline{y}, X) = \frac{|X^t X|^{1/2}}{(\sqrt{2\pi})^k \sigma^k} \exp \left\{ -\frac{1}{2\sigma^2} [(\underline{\theta} - \hat{\underline{\theta}})^t X^t X (\underline{\theta} - \hat{\underline{\theta}})] \right\} \quad -\infty < \theta_i < \infty, \quad \forall i = 1, \dots, k \quad (1.21)$$

de la expresión (1.21), se puede notar que lo único relevante acerca de  $\underline{\theta}$  en la distribución, fue la información presentada en la muestra. El vector de medias se distribuye como una Normal Multivariada con media  $\hat{\underline{\theta}}$  y varianza  $\sigma^2 (X^t X)^{-1}$ ; es decir

$$\underline{\theta} \sim N_k(\hat{\underline{\theta}}, \sigma^2 (X^t X)^{-1}). \quad (1.22)$$

Para propósitos de inferencia o decisión acerca del parámetro  $\underline{\theta}$  pasar a la sección 1.1.1 o 1.1.2, respectivamente

### 1.2.2 $\underline{\theta}$ Y $\sigma^2$ DESCONOCIDOS.

Como ambos parámetros son desconocidos la distribución final consistirá de

$$p(\underline{\theta}, \sigma^2 | \underline{y}, X) \propto p(\underline{\theta}, \sigma^2) p(\underline{\theta}, \sigma^2 | \underline{y}, X), \quad (1.23)$$

y usando la dependencia de las estadísticas  $\hat{\underline{\theta}}$  y  $S^2$ , dados  $\underline{\theta}$  y  $\sigma^2$ , se puede escribir (1.23) de otra forma, en términos de las estadísticas suficientes

$$p(\underline{\theta}, \sigma^2 | \underline{y}, X) \propto p(\underline{\theta}, \sigma^2) p(\hat{\underline{\theta}} | \underline{\theta}, \sigma^2) p(S^2 | \sigma^2). \quad (1.24)$$

Por lo que ahora utilizando la distribución inicial de referen-

$$p(\underline{\theta} | y, x) = \frac{|x^t x|^{-k/2}}{(\sqrt{2\pi})^k \sigma^k} \exp \left\{ -\frac{1}{2\sigma^2} [y - \underline{\theta}^t x]^2 \right\} \quad -\infty < \theta_j < \infty, \quad k = 1, \dots, K \quad (1.21)$$

de la expresión (1.21), se puede notar que lo único relevante acerca de  $\underline{\theta}$  en la distribución, fue la información presentada en la muestra. El vector de medias se distribuye como una Normal Multivariada con media  $\underline{\theta}$  y varianza  $\sigma^2 (x^t x)^{-1}$ ; es decir

$$\underline{\theta} \sim N_k(\underline{\theta}, \sigma^2 (x^t x)^{-1}). \quad (1.22)$$

Para propósitos de inferencia o decisión acerca del parámetro pasar a la sección 1.1.1 o 1.1.2, respectivamente

### 1.2.2 $\underline{\theta}$ Y $\sigma^2$ DESCONOCIDOS.

Como ambos parámetros son desconocidos la distribución final consistirá de

$$p(\underline{\theta}, \sigma^2 | y, x) \propto p(\underline{\theta}, \sigma^2) p(\underline{\theta}, \sigma^2 | y, x), \quad (1.23)$$

y usando la dependencia de las estadísticas  $\hat{\underline{\theta}}$  y  $\hat{\sigma}^2$ , dados  $\underline{\theta}$  y  $\sigma^2$ , se puede escribir (1.23) de otra forma, en términos de las estadísticas suficientes

$$p(\underline{\theta}, \sigma^2 | y, x) \propto p(\underline{\theta}, \sigma^2) p(\hat{\underline{\theta}} | \underline{\theta}, \sigma^2) p(\hat{\sigma}^2 | \sigma^2). \quad (1.24)$$

Por lo que ahora utilizando la distribución inicial de referen-

cia para los parámetros desconocidos en (1.17), tenemos que también se puede escribir como

$$p(\theta, \sigma^2) \propto \sigma^{-4} \quad \begin{matrix} -\infty < \theta < \infty \\ 0 < \sigma < \infty \end{matrix} \quad (1.25)$$

Box y Tiao (1971, pp.62) encuentran una expresión para  $p(\sigma^2/\sigma^2)$  y obtienen la constante de proporcionalidad para expresar (1.23), es decir, reexpresar a  $p(\theta, \sigma^2/\underline{y}, X)$  de la siguiente manera

$$p(\theta, \sigma^2/\underline{y}, X) = p(\sigma^2/S^2) p(\theta/\hat{\theta}, S^2), \quad (1.26)$$

además, se determinaron cada uno de los términos en (1.26); la densidad  $p(\sigma^2/S^2)$  le corresponde una ji-cuadrada invertida (Box y Tiao, pp.63) y la densidad  $p(\theta/\hat{\theta}, S^2)$ , se notó que es la misma que en  $p(\theta/\underline{y}, X)$  (1.21). Entonces, se puede expresar la posterior de ambos parámetros desconocidos mediante

$$p(\theta, \sigma^2/\underline{y}, X) = \frac{1}{\Gamma(\frac{n-1}{2})} \left(\frac{n-1}{2\sigma^2}\right)^{\frac{n-1}{2}} (S^2)^{\frac{n-1}{2}-1} e^{-\frac{(n-1)S^2}{2\sigma^2}} \frac{|X^t X|}{(\sqrt{2\pi}\sigma)^k} e^{-\frac{1}{2\sigma^2}(\theta - \hat{\theta})^t X^t X (\theta - \hat{\theta})} \quad (1.27)$$

y reagrupando (1.27), se tiene

$$p(\theta, \sigma^2/\underline{y}, X) = \frac{1}{\Gamma(\frac{n-1}{2})} \left(\frac{n-1}{2\sigma^2}\right)^{\frac{n-1}{2}} (S^2)^{\frac{n-1}{2}-1} \frac{|X^t X|}{(\sqrt{2\pi}\sigma)^k} e^{-\frac{(n-1)S^2}{2\sigma^2} + (\theta - \hat{\theta})^t X^t X (\theta - \hat{\theta})} \quad (1.28)$$

Se concluye que la densidad de  $\theta$ , posee a  $\hat{\theta}$  el llamado estimador de "Mínimos Cuadrados". Este valor minimiza la forma cuadrática

$Q(\theta) = (\underline{y} - X\hat{\theta})^t X^t X (\underline{y} - X\hat{\theta})$ , y se expresa como en (1.5) y (1.6).

De necesitar alguna decisión o inferencia iniciar con la sección 1.1.2 ó 1.1.1, respectivamente y seguir de nuevo aquí (en caso de inferencia).

Box y Tiao (1971, pp. 125 - 127) hacen una justificación bayesiana del empleo de la tabla de Análisis de Varianza, que es comunmente utilizada en la Estadística Clásica, por medio de los siguientes argumentos; la forma cuadrática de  $p(\theta/y)$  definida por  $KS^2$  se distribuye como una  $F$  con ciertos parámetros conocidos, (esto es,  $Q(\theta)/S^2K \sim F_{(k,v)}$ ). Para saber si un punto  $\theta_0$  se encuentra en la región de alta densidad  $I_A$  es suficiente y necesario que  $(\theta_0 - \hat{\theta})' X' X (\theta_0 - \hat{\theta}) < KS^2 F_{(k,v,\alpha)}$  (Este resultado está relacionado con la expresión (1.13)). O equivalentemente calcular la cantidad

$$P_r \left\{ F_{(k,v)} < \frac{(\theta_0 - \hat{\theta})' X' X (\theta_0 - \hat{\theta})}{KS^2} \right\} \quad (1.29)$$

donde  $F_{(k,v)}$  es una variable con  $(k,v)$  grados de libertad, da la probabilidad contenida de la región de alta densidad que incluye a  $\theta_0$ .

La región de alta densidad de contenido  $1-\alpha$  es numéricamente idéntica a la región de confianza de tamaño mínimo, y el complemento de la probabilidad en (1.29) da el nivel de significancia asociado con la hipótesis nula  $H_0: \theta = \theta_0$  contra la alternativa  $H_a: \theta \neq \theta_0$ .

Así clásicamente en la tabla de Análisis de Varianza, se determina si un parámetro está incluido en la región de alta densidad por medio del cálculo de los Cuadrados Medios y el cociente de la  $F_{(k,v,\alpha)}$  de tablas.

A continuación se presenta la tabla de Análisis de Varianza

FUENTE	SUMA DE CUADRADOS	GRADOS DE LIBERTAD	CUADRADOS MEDIOS	COCIENTE DE CUADRADOS MEDIOS.
Discrepancia de parámetros	$(\theta_0 - \hat{\theta})^t X (\theta_0 - \hat{\theta}) = A$	$K$	$A/K = C$	$\%D$
Residual	$(y - \hat{y})^t (y - \hat{y}) = B$	$n - K$	$B/(n - K) = D$	
Total	$(y - X\theta_0)^t (y - X\theta_0)$	$n$	$1$	

En ocasiones, se desea inferir sobre combinaciones de los parámetros y para tales casos, se determina una función  $\Phi$ , con  $m$  parámetros definidos de tal manera que  $m \leq K$  y

$$\Phi = (\phi_1, \dots, \phi_m)^t, \quad \phi_i = f_i(\theta) \quad i = \overline{1, m}, \quad (1.30)$$

donde  $\Phi$  es una transformación lineal de  $\theta$  ( $\theta \in R^K$ ), mediante la obtención de densidad correspondiente, por lo que pueden realizarse las inferencias deseadas.

Las comparaciones de parámetros que se pueden efectuar es mediante las regiones de alta densidad y surgen a partir de problemas estadísticos importantes como: comparaciones de medias o de varianzas en  $K$  distribuciones normales.



## 2. GENERALIZACIÓN DEL MODELO LINEAL.

Una generalización del modelo lineal (Box y Tiao, pp. 176-178) consiste en la aplicación del modelo de regresión lineal a observaciones que no tienen una distribución normal en los errores, pero que pertenecen a la familia exponencial potencia.

El modelo lineal general es entonces

$$y = X\theta + \epsilon, \quad (1.31)$$

donde:

$y$  es un vector de observaciones,  $X$  una matriz de rango completo con elementos fijos,  $\theta$  vector de coeficientes de regresión desconocidos y  $\epsilon$  vector de errores aleatorios.

Suponiendo que los elementos de los errores son independientes y que tienen una distribución exponencial potencia definida por

$$p(\epsilon | \theta, \beta) = w(\beta) \sigma^{-1} e^{-c(\beta) \left(\frac{\epsilon}{\sigma}\right)^{\frac{1}{\beta} + 1}} \quad -\infty < \epsilon < \infty \quad (1.32)$$

con

$$w(\beta) = \frac{\left\{ \Gamma\left[\frac{3}{2}(1+\beta)\right] \right\}^{1/2}}{(1+\beta) \left\{ \Gamma\left[\frac{1}{2}(1+\beta)\right] \right\}^{3/2}} \quad 0 > 0, -\infty < \theta < \infty$$

$$-1 < \beta < 1$$

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

Como última suposición se pide que el modelo no tenga problemas

de multicolinealidad.

A los parámetros  $\sigma$  y  $\beta$  son conocidos como la desviación estándar y medida de la curtosis, respectivamente. El último parámetro  $\theta$ , es una medida de la no normalidad en la distribución origen. Por ejemplo, si  $\theta = 0$  resulta ser la conocida normal y con  $\theta = 1$ , la doble exponencial.

Las complicaciones en las expresiones de la función de verosimilitud son consecuencias del manejo de tres parámetros  $\{\theta, \sigma, \beta\}$  y su representación es como sigue:

$$l(\theta, \sigma, \beta | y) \propto (\omega(\beta))^n \delta^{-n} \mathcal{L}^{-c(\beta)} \prod_{i=1}^n \left| \frac{y_i - \lambda_{(i)}(\theta)}{\sigma} \right|^{2/\theta}, \quad (1.33)$$

con

$$-\infty < \theta < \infty, \quad \sigma > 0, \quad -1 < \beta < 1$$

$\lambda_{(i)}^t$  es el  $i$ -ésimo renglón de  $X$ .

Suponiendo independencia entre los parámetros se puede expresar la densidad inicial conjunta como

$$p(\theta, \sigma, \beta) = p(\beta) p(\theta, \sigma), \quad (1.34)$$

que es una manera de manejar separadamente el parámetro de la curtosis y los otros parámetros  $\{\theta, \sigma\}$ . Utilizando en  $p(\theta, \sigma)$  una inicial de referencia, la densidad resulta como:

$$p(\theta, \sigma, \beta) \propto (\omega(\beta))^n \delta^{-(n+\theta)} \mathcal{L}^{-c(\beta)} \prod_{i=1}^n \left| \frac{y_i - \lambda_{(i)}(\theta)}{\sigma} \right|^{2/\theta} p(\theta). \quad (1.35)$$

Integrando con respecto a  $\theta$ , la expresión (1.35), Sox y Tiao la expresa en términos de

$$p(\underline{\theta}, \beta | \underline{y}) = p(\underline{\theta} | \beta, \underline{y}) p(\theta | \underline{y}), \quad (1.36)$$

y tomando únicamente la distribución condicional  $p(\underline{\theta} | \beta, \underline{y})$  en (1.36) se obtiene

$$p(\underline{\theta} | \beta, \underline{y}) = \left[ \int_{\mathcal{R}} [M(\theta)]^{-\frac{n}{2}(1+\beta)} d\theta \right]^{-1} [M(\theta)]^{-\frac{n}{2}(1+\beta)}, \quad \begin{matrix} -\infty < \theta_j < \infty \\ j = 1, \dots, K \end{matrix} \quad (1.37)$$

donde

$$M(\theta) = \prod_{i=1}^n |y_i - x_{(i)}^t(\theta)|^{2/1+\beta}$$

El estudio de la distribución  $p(\underline{\theta}, \beta, \underline{y})$  como una función de  $\beta$ , puede determinar el tipo de desviación de la normalidad. Por ejemplo, en el caso de tener  $\beta = 0$ ,  $M(\theta) = (\underline{y} - X\theta)^t (\underline{y} - X\theta)$  y la distribución en  $p(\underline{\theta} | \beta, \underline{y})$  es una  $t_K(\underline{\theta}, S^2(X^t X)^{-1}, \nu)$ . Para otros casos cuando  $\theta$  es pequeña ( $K=2$  ó  $K=3$ ), los contornos de  $M(\theta)$  únicamente se pueden graficar.

Como se ve, en el caso de no contar con normalidad en los errores, puede resultar muy difícil el manejo práctico de la distribución final, por las expresiones tan grandes. Sería conveniente efectuar pruebas de hipótesis en  $p(\theta | \underline{y})$  (1.36) para inferir si hay o no separación de normalidad (es decir,  $H_0: \beta = 0$  vs  $H_a: \beta \neq 0$ ) y si no hay mucho problema de la normalidad se podrá hacer una transformación o cambiar el modelo a otro más adecuado, pero si la evidencia apunta a la no normalidad se tendrá que trabajar con el modelo lineal general.



**METODOS DE  
SELECCION  
DE VARIABLES**

## CAPITULO I I

### MÉTODOS DE SELECCION DE VARIABLES.

La idea de este capítulo es la de exhibir algunos de los métodos que existen en la Estadística Bayesiana y cuyo objetivo es la selección de variables en problemas de regresión.

El problema relacionado con Selección de Variables resulta ser un caso particular de la selección de modelos, es decir, cuando en la estadística se presenta el problema de escoger un conjunto de variables explicativas en el Análisis de Regresión, se está seleccionando a la vez el modelo que está constituido por aquellas variables.

Debido a las dificultades que produce la Multicolinealidad en los modelos (Draper y Smith(1981, pp.258) y Chatterjee(1977, pp.143-192)) de regresión para inferir y predecir, es conveniente hacer una división en la discusión de los conjuntos de variables que presentan este problema, para aplicarles "Un trato especial" previo a los procedimientos de Selección de Variables.

Los conjuntos de variables que no tengan problemas de Multicolinealidad serán directamente sometidos a los métodos de selección de variables que se expondrán más adelante.

Existen diversas técnicas que indican y permiten corregir la Mul

colinealidad, aquí únicamente se comentará el método de componentes principales (para una discusión más detallada ver Draper y Smith o Chatterjee).

Habiendo eliminado de alguna manera el problema anterior, es necesario validar si en realidad el modelo con el conjunto de todas las variables (modelo completo) explica o predice la variable dependiente, de no ser así, no tendrá sentido efectuar una selección de variables que no satisfacen ésta propiedad, pues cualquier subconjunto de ellas lo hará.

Se comenzará suponiendo que el modelo completo está formado por  $p$  parámetros (o variables explicativas). De aquí, que el número total de posibles modelos a analizar con respecto al modelo completo son  $2^p - 1$ . Para lo que si no es muy grande esto crece mucho; inclusive para valores relativamente pequeños de  $p$ , esto se agranda, por ejemplo si  $p$  es igual a 10 resultan 1022 modelos diferentes para analizar.

Para algunos estadísticos la elección del "Mejor modelo" es consecuencia de la bondad (aproximada) de su predicción o explicación; de lo fácil y sencillo del manejo de datos, así como también de su costo de captación de datos. Aunque estas características en general pueden cambiar dependiendo del problema particular y de sus intereses.

Los procedimientos en la literatura bayesiana que se encontraron se clasificaron en dos secciones: 2.2.1 la sección de inferencia, donde se agrupa a los procedimientos que no tienen objetivos predictivos implícitamente, y la segunda sección, 2.2.2 donde se involu-

eran a los métodos que especifican tener metas predictivas. Esta última sección se divide en dos subsecciones: 2.2.2.1 artículos que no se desarrollan en base a la Teoría de Decisiones y 2.2.2.2. artículos involucrados en la teoría de decisiones (porque involucran una función de pérdida).

## 2.1 MULTICOLINEALIDAD.

El surgimiento del problema de Multicolinealidad o también llamado datos colineales se debe principalmente a 1) un inadecuado diseño de muestreo, 2) una sobredefinición del modelo, pues existen más variables que observaciones y 3) las características inherentes al procedo de investigación, que ocasionan interrelaciones entre las variables explicativas.

Los síntomas para sospechar Multicolinealidad principalmente son: 1) grandes cambios en los coeficientes estimados cuando una variable es sumada o eliminada, 2) grandes cambios en los coeficientes cuando se excluye o incluye en el análisis un dato, 3) resultados no esperados tanto en los signos como en los coeficientes estimados, y 4) los coeficientes de las variables que se esperaría fueran importantes tienen grandes errores estándar.

Otra indicación para detectar Multicolinealidad se presenta en la matriz de correlación de las variables explicativas, ya que aparecen fuera de la diagonal correlaciones proximas a 1, se dirá que hay Multicolinealidad. Esta indicación no es muy buena, ya que puede existir Multicolinealidad producida por alguna agrupación de las variables explicativas y no indicarse en la matriz de correlación.

El procedimiento de componentes principales es una herramienta estadística del Análisis Multivariado (Chatfield y Collins (1980, pp. 57-62)) que consiste en construir un conjunto de nuevas variables explicativas ortogonales ( $Z_1, Z_2, \dots, Z_n$ ), que son formadas con base en combinaciones lineales de las variables explicativas originales ( $X_1, \dots, X_n$ ).



Componentes principales se aplica tanto sobre la matriz  $P$  de correlaciones como sobre la matriz de varianzas-covarianzas,  $D$  (se obtienen resultados distintos al usar estas matrices). La matriz de varianzas y covarianzas de las nuevas variables,  $A$ , es una matriz diagonal cuyos elementos son iguales a las raíces características  $\lambda_i$  ( $i = \overline{1, n}$ ) de la matriz  $P$  o  $D$  (dependiendo de la matriz que se usó). Estos valores propios son las varianzas respectivas de las nuevas componentes. Por construcción la suma de las varianzas de las variables originales y la suma de las componentes principales son iguales ( $\sum_{i=1}^n \text{Var}(Z_i) = \sum_{i=1}^n \lambda_i = \sum_{i=1}^n \text{Var}(X_i)$ ).

En el caso de utilizar inicialmente la matriz  $P$  al aplicar componentes principales, la matriz  $A$  será tal, que en su diagonal puede haber solo unos. Si esto ocurre resulta que las variables  $X_i$  son ortogonales. Pero si este no es el caso, hay al menos una  $\lambda_i = 0$  o muy próxima a cero, entonces se dice que existe Multicolinealidad en las variables explicativas originales.

Cuando los resultados apoyen la existencia de datos colineales hay ciertos métodos en estadística que permiten elegir en su conjunto de componentes principales. Uno de ellos, consiste en elegir las primeras componentes principales<sup>(1)</sup>, digamos  $m$ , de tal manera que en conjunto expliquen cierto porcentaje de variación total en los datos originales. Por ejemplo, se aceptan las  $m$  primeras  $\lambda_i$  si  $\frac{\sum_{i=1}^m \lambda_i}{\sum_{i=1}^n \lambda_i}$  rebasa cierto porcentaje establecido. Otros criterios analizan a la matriz  $A$ , de la siguiente forma: si  $\lambda$  es aproximadamente cero (generalmente se empiezan analizando la última  $\lambda$ ,  $\lambda_n$ ) entonces la correspondiente componente principal es aproximadamente constante, de donde se concluye que no se está aportando ninguna variabilidad y esto permite estudiar a la Multicolinealidad. De esta manera, se puede expresar a una de las variables originales como combinación lineal de las otras.

(1) Esto se debe a la construcción de componentes principales,  $\lambda_1 \geq \lambda_2 \geq \dots, \lambda_n$ .

El paso siguiente consiste en pedir información (extra-matemática) al interesado sobre la situación inherente del problema, para que apoyados en las información<sup>(1)</sup> evalúen la posibilidad de eliminar alguna combinación lineal de las otras. Todo esto con el objetivo de reducir el número de variables originales al rango de la matriz  $\Sigma(SP)$ . En el caso de no poder eliminar alguna variable original, se expresa, cada una de las componentes principales diferentes de cero, en términos de las variables  $X$ 's para plantear el modelo de regresión obtenido. Por ejemplo, aplicando componentes principales se obtienen diferentes  $\lambda$ 's, de las cuales, hay  $\kappa$  diferentes de cero ( $1 \leq \kappa < n$ ) y las restantes  $n - \kappa$ , se determinaron cero. En base a lo anterior, se expresa al modelo de regresión que explica  $\delta$  predice a  $Y$ , en términos de las  $\kappa$  componentes principales diferentes de cero, es decir, a las componentes que se determinaron como cero, se estandarizan con respecto a su media y se les suman a término constante del modelo, lo cual permite plantear al modelo en términos de las  $\kappa$  componentes principales. Substituyendo en cada una de las  $i$ 's por las combinaciones respectivas de las  $X$ 's, queda resuelto el problema de Multicolinealidad sobre un subespacio (o dimensión menor) del problema original.

Un aspecto importante sobre el procedimiento de componentes principales, es en el siguiente sentido, no se asegura que los modelos obtenidos por este procedimiento sean "buenos" para los objetivos específicos en que fueron creados, solo se sabe que estos nuevos datos no presentan el problema de multicolinealidad y que sus componentes son combinaciones lineales de los datos originales. Por lo que estos modelos pueden ser malos si se trata de explicar a la variable dependiente.

(1) Esta situación va ir de acuerdo con las correlaciones de las variables explicativas.

## 2.2 MÉTODOS DE SELECCIÓN DE MODELOS,

Antes de iniciar la presentación de los métodos de selección de variables se establecerá una terminología general.

Tomando en cuenta el anterior capítulo, se recordará que el problema de regresión involucra en esencia tres términos: 1) los parámetros  $\underline{\theta}$  y  $\omega$ , donde  $\underline{\theta}$  es un vector en  $\mathbb{R}^p$  ( $\underline{\theta} = (\theta_1, \dots, \theta_p)^t$ ) y  $\omega$  un parámetro real positivo, 2) un vector de observaciones comúnmente conocido como  $\underline{y}$  de tamaño  $1 \times p$  ( $\underline{y} = (y_1, \dots, y_n)^t$ ) y por último 3) una matriz de datos conocidos de tamaño  $n \times p$ , que se puede reexpresar como  $p$  vectores columna de tamaño  $1 \times p$  ( $X = (\underline{x}_1, \underline{x}_2, \dots, \underline{x}_p)^t$ ).

Suponiendo normalidad en los errores, se tiene que la verosimilitud de la muestra de los datos, es una distribución normal  $n$ -variada con media  $X\underline{\theta}$  y precisión  $\omega I_n$  ( $p(\underline{y}/\underline{\theta}, \omega, X) \sim N_n(X\underline{\theta}, \omega I_n)$ ).

Por comodidad se tomará la siguiente notación para hablar de modelos que contengan  $j$  variables explicativas con  $j = \overline{1, p}$ . Se le asignará el nombre de modelo completo a  $M_c: f(x_1, \dots, x_p)$ , a todos los demás subconjuntos de variables  $Z_i$  que tengan menos de  $p$  variables explicativas se le asociará un modelo  $M_i$  que representa un modelo de regresión sobre las variables explicativas en  $Z_i$  con  $i = \overline{1, 2^p - 2}$ .

Luego entonces, el modelo  $i$ -ésimo tendrá la siguiente forma

$$M_i: f(Z_i) = Z_i \underline{\theta}_i + \varepsilon$$

donde  $z_i$  representa las variables explicativas incluidas en el modelo  $M_i$  y  $\theta_i$  representa el vector de parámetros correspondientes.

Cabe aclarar un punto importante acerca de los procedimientos de que se presentaran. Los eventos de interés relacionados con los modelos no han sido establecidos en todos los métodos de selección, ya que algunos de éstos (por ejemplo: Zellner y Atkinson) no especifican a que evento relacionado con el modelo se refieren. No se sabe si hablan de la probabilidad de que el modelo sea el verdadero o a la probabilidad de que el modelo contenga un conjunto determinado de variables o tal vez se ha dejado abierta la especificación de este evento para que dependiendo del caso y las necesidades del interesado se establezcan. Pero como se verá más adelante, existen otros procedimientos de selección (por ejemplo: Davis o Geisser y Eddy) en donde sí establecen el evento asociado a los modelos, es la probabilidad de que el modelo sea correcto.

### 2.2.1 INFERENCIA.

Los procedimientos que se presentaran aquí, están desarrollados por Zellner, Atkinson, Davis, Broemeling y Perichi.

#### PROBABILIDADES POSTERIORES.

Este procedimiento estriba en elegir un modelo con mayor probabilidad posterior, o sea, el método consiste en la comparación de hipótesis, en donde la información inicial del modelo puede ser involucrada en las probabilidades posteriores (Zellner, 1971, pp. 291-298).

Las hipótesis involucradas en este método se relacionan con los modelos de regresión  $M_i$ , donde estos pueden ser de  $i = \overline{1, m}$  y  $m = 2^k - 2$ , considerando una posible eliminación de los modelos por razones no estadísticas.

Otra característica de estos procedimientos, es que presentan cualidades de exclusividad y exhaustividad, esto es, deben incluir a todos los posibles modelos aceptables y entre ellos, debe de estar incluido el "Modelo Satisfactorio", o como algunos autores lo llaman "Modelo Verdadero".

La probabilidad posterior del modelo  $i$  ( $i = \overline{1, m}$ ) con respecto a los datos, es decir  $p(M_i / X, \underline{\theta}, \underline{Z}_i, \underline{Y})$  se encuentra por medio del producto de la probabilidad asignada al modelo  $i$  por la verosimilitud del modelo  $i$  fijando a los datos, entre todos los modelos con iguales restricciones, es decir,

$$P(M_i | \text{datos}) = P(M_i | X, \theta, \omega, z_i, Y) = \frac{P(M_i) L(M_i | \text{datos})}{\sum_{i=1}^2 P(M_i) L(M_i | \text{datos})} \quad (2.2)$$

Cada verosimilitud de  $M_i$  con respecto a los datos, es el resultado del producto de  $n$  funciones de distribución,  $p(y_j | M_i, z_j, x_j)$  que involucra a la vez las funciones de densidad,  $f_j(y_j | z_j, \theta_j, \omega)$ , para  $Y$ , además la función de densidad inicial  $f(\theta_j; \omega)$ , se expresa como

$$L(M_i | \text{datos}) = \prod_{j=1}^n p(y_j | M_i, z_j) = \prod_{j=1}^n \int_{\mathbb{R}^k} f_j(y_j | z_j, \theta_j, \omega) f(\theta_j; \omega) d\theta_j d\omega.$$

Uno de los pioneros que trabajó este procedimiento de probabilidades posteriores fue Zellner (1971, pp. 291-298) que desarrolló su trabajo en el caso más simple.

Zellner consideró únicamente la existencia de dos modelos  $(M_1, M_2)$  que explican la variación de la variable dependiente  $y$  y donde sólo uno de ellos es el verdadero (aquí,  $M_1: H_0$  y  $M_2: H_0^c$ ).

El vector de observaciones  $y^t = (y_1, \dots, y_n)$  está dado en la forma general (1.1), que en esta situación se expresa como:

$$M_1: y = X_1 \beta_1 + u_1$$

$$M_2: y = X_2 \beta_2 + u_2,$$

donde  $X_1$  y  $X_2$  son matrices de  $n \times k$  con cantidades conocidas y cada una

con rango  $K$ ,  $A_1$  y  $A_2$  vectores de coeficientes  $(K \times 1)$  con elementos no comunes,  $u_1$  y  $u_2$  vectores de errores  $(n \times 1)$ , por último, para el modelo verdadero  $M_1$  o  $M_2$  se asume que existe independencia en los errores y se distribuye como una normal con media cero y varianzas  $\sigma_i^2$ ,  $i = 1, 2$ .

Se hacen ciertas suposiciones distribucionales para los parámetros que consisten en  $p(\beta_i, \sigma_i) = p(\beta_i / \sigma_i) p(\sigma_i)$ ; se pide que  $p(\beta_i, \sigma_i)$  se distribuya como una normal con media  $\bar{\beta}_i$  y varianzas  $\sigma_i^2 C_i^{-1}$  ( $C_i$  es una matriz simétrica definida positivamente de  $K \times K$ ) y  $p(\sigma_i)$  se distribuye como una Gamma invertida (Anexo 1) con parámetros  $q_i$  y  $S_i^{-1}$  que son asignados por el investigador ( $q_i > 0$  y  $\bar{S}_i, \alpha_i = 1, 2$ ).

Se concluye llegando al "Modio" posterior dado por

$$K_{1,2} = \frac{P(M_1) \left[ \frac{|C_1|}{|A_1|} \right]^{1/2} \left( \frac{\sigma_1}{\sigma_2} \right)^{-n/2} \frac{\bar{S}_1^{-1} / \sigma_1 + q_{1,n} \left( \frac{\bar{S}_1^{-1}}{\sigma_1} \right)}{P(M_2) \left[ \frac{|C_2|}{|A_2|} \right]^{1/2} \left( \frac{\sigma_2}{\sigma_1} \right)^{-n/2} \frac{\bar{S}_2^{-1} / \sigma_2 + q_{2,n} \left( \frac{\bar{S}_2^{-1}}{\sigma_2} \right)}{2.3}$$

con

$$J_i = \frac{1}{n} \left[ v S_i^{-1} + (\hat{\beta}_i - \bar{\beta}_i)' C_i (\hat{\beta}_i - \bar{\beta}_i) + (\hat{\sigma}_i - \bar{\sigma}_i)' X' X (\hat{\sigma}_i - \bar{\sigma}_i) \right], \quad i = 1, 2$$

$$A_i = C_i + X_i' X_i, \quad \bar{\beta}_i = A_i^{-1} (C_i \bar{\beta}_i + X_i' X_i \hat{\beta}_i)$$

$$\hat{\beta}_i = (X_i' X_i)^{-1} X_i' y, \quad \bar{\beta}_i \text{ es el vector de medias inicial para } \beta_i$$

$$v S_i^{-1} = (y - X_i \hat{\beta}_i)' (y - X_i \hat{\beta}_i), \quad v = n - K$$

$$q_{i,n} \left( \frac{S_i^{-1}}{\sigma_i} \right)$$

$k_{12}(\bar{S}_1^2 / S_1)$  denota la ordenada de la distribución  $F$  con  $q_1$  y  $n$  grados de libertad.

Esta expresión resulta ser un caso particular del cociente (2.2) de probabilidades posteriores.

La expresión (2.3) cambia bajo algunas de estas condiciones; si  $n$  es grande se asume  $\frac{|C_1|/|A_1|}{|C_2|/|A_2|} \rightarrow 1$ , y  $k_{12} = \left(\frac{S_1^2}{S_2^2}\right)^{-\frac{n+4}{2}} \left(\frac{q_1 S_1^2}{2S_2^2} + \frac{q_2 S_2^2}{2S_1^2}\right)$ , pero si además se cuenta con información inicial "difusa"  $k_{12}$ , se reduce a

$$k_{12} = \left(\frac{S_1^2}{S_2^2}\right)^{-\frac{n}{2}} \quad (2.4)$$

La importancia de estos resultados radica, en que si se cuenta con una función de pérdida simétrica (ver anexo) se elegirá aquel modelo con mayor probabilidad posterior y bajo las condiciones de Zellner, es elegir el modelo con mayor coeficiente de determinación o que es lo mismo que elegir el modelo con menor  $S^2$ .

Como último comentario al trabajo de Zellner, este no fue ubicado en los procedimientos de Teoría de Decisiones porque la función de pérdida que se menciona es muy particular y probablemente no muy común. Esto se debe a que únicamente se tiene dos estados de la naturaleza ( $\theta_1, \theta_2$ ) y ambos tienen asociados la misma pérdida en el caso de haber tomado malas decisiones ( $L(d_1, \theta_2) = L(d_2, \theta_1) = c$ ) y pérdida cero si la decisión es correcta. Además, si se toman ahora tres estados de la naturaleza ( $\theta_1, \theta_2, \theta_3$ ) con tres decisiones ( $d_1, d_2, d_3$ ) se concluye que el maximizar la probabilidad no coincide siempre con la decisión de minimizar la pérdida esperada. Por ejemplo, si se tiene la siguiente función de pérdidas, en donde la decisión correc



ta trae una pérdida cero; es decir,  $L(d_1, \theta_1) = L(d_2, \theta_2) = L(d_3, \theta_3) = 0$ , la pérdida asociada en el caso de haber tomado una mala decisión es la misma en todos los casos  $L(d_1, \theta_1) = L(d_1, \theta_2) = L(d_2, \theta_1) = L(d_2, \theta_3) = L(d_3, \theta_1) = L(d_3, \theta_2) = c$ , y la probabilidad asignada a los estados de la naturaleza es la siguiente  $P(\theta_1) = .4$  y  $P(\theta_2) = .3 = P(\theta_3)$ . En el caso de maximizar probabilidades se elige como óptimo a  $d_1$  y en el caso de minimizar la pérdida esperada se elige como decisión óptima también a  $d_1$ . Por lo que se concluye que en este caso coinciden ambos criterios. Otro ejemplo, basado en el anterior enunciado se mantienen fijas las probabilidades asignadas a los estados de la naturaleza y se establece la función de pérdidas, como a continuación se especifica:

$$L(d_1, \theta_1) = 0, L(d_1, \theta_2 \cup \theta_3) = .6c, L(d_2, \theta_1) = L(d_3, \theta_1) = .4c,$$

$$L(d_2, \theta_2 \cup \theta_3) = L(d_3, \theta_2 \cup \theta_3) = .3c,$$

por lo cual con el criterio de maximizar las probabilidades se elige a  $d_2$  ó  $d_3$ , y en minimizar la pérdida esperada se elige a  $d_1$ , por lo tanto, hay contradicción y ésta se debe a que el criterio de maximizar probabilidades posteriores viola los axiomas de coherencia, por lo tanto es una mala decisión el usarlo.

Otro trabajo basado en probabilidades posteriores, es el de Atkinson (1978), en donde considera el caso general con  $m$  modelos con las características ya antes mencionadas, teniendo entonces modelos en competencia, es decir, todas son modelos que predicen o explican y sólo una de estos puede ser elegido. Una condición que se pedirá es que todos los modelos contengan el mismo número de parámetros, porque de lo contrario acarrearla inferencias arbi-

trarias como por ejemplo, el favorecer a modelos con pocos parámetros.

Otra condición que se pide es que la suma sobre las probabilidades de los modelos condicionados a los datos sume uno.

Atkinson indica que cuando  $\sigma^2$  es desconocida en los modelos de Regresión, no existe ninguna dificultad en este procedimiento, por lo que resulta una expresión análoga a (2.3) y que únicamente se modificará su distribución por una *t*-student.

Mediante ejemplos de dos modelos en diferentes circunstancias se exponen las fallas y alcances que se presentan a utilizar este procedimiento de probabilidades posteriores, las cuales son: 1) que si los modelos tienen igual número de parámetros y la información inicial es no informativa, entonces en este procedimiento se favorece al modelo incorrecto, 2) para el caso de tener modelos incorrectos, el procedimiento elige de todos los modelos, un modelo cualquiera. De aquí, que se requiera checar los modelos adecuados por medio de un estimador de la varianza; para revisar posible carencia de ajuste, 3) si se trabaja con modelo "separados" (modelos que no están anidados<sup>(1)</sup>, y tienen variables explicativas distintas), entonces las probabilidades posteriores tienden a favorecer asintóticamente al modelo "correcto" ó al modelo con menor número de parámetros, y 4) con modelos anidados correctos se obtienen resultados que apoyan al modelo con menor número de parámetros.

Atkinson hace la suposición de la existencia de un modelo verda-

(1) Se dice que el modelo A está contenido en el modelo B, si todas las variables explicativas de A son variables explicativas en B.

dero, que en la práctica no necesariamente es único, como en el caso de regresión con modelos anidados. Aclarando esto, en caso de tener un problema de Selección de Variables en Regresión para algunos subconjuntos de variables se tienen modelos anidados, por lo que al asignarle una probabilidad a uno de estos modelos (si el evento relacionado con el modelo es obtener el mejor modelo), se le tendría que asignar una probabilidad mayor o igual a todos aquellos modelos que tengan esas variables explicativas. Esto ocasionaría que el modelo con mayor probabilidad resulte ser el modelo completo, ocasionando que el criterio de probabilidades posteriores siempre elija el modelo completo.

También existe otro problema que no es propio de este procedimiento de selección, pues se presenta en todos los métodos clásicos o bayesianos. En este caso se ve muy claro cuando el investigador asigna probabilidades a los modelos considerados y no ha podido establecer a "todos" los modelos posibles. Como en este método no es factible tener la opción "otros modelos" resulta que "otros modelos" tienen probabilidades cero, y puede suceder que entre los no considerados existan algunos mejores modelos. Aunque siempre exista un modelo que no se tome en cuenta, se espera que a pesar de ello, se acerquen los resultados a los "verdaderas" soluciones.

Si siguiendo con este procedimiento de Probabilidades Posteriores se tiene el trabajo realizado por Davis (1979), que se proporciona una cota para el error obtenido al calcular la probabilidad productiva cuando una densidad inicial aproximada<sup>(1)</sup> es usada. Es decir, partiendo del principio de estimación estable<sup>(2)</sup>, Davis encuentra

- (1) Algunos estadísticos creen en la existencia de densidades iniciales verdaderas y por lo tanto en densidades posteriores verdaderas. Con esta idea se establecen posibles aproximaciones a ellas.
- (2) Este principio supone la existencia de un conjunto sobre el cual para una muestra lo suficientemente grande, la distribución inicial es dominada por los datos. Trabajo realizado por Lindeman, Savage y Edwards, en 1963.

una cota para el porcentaje de error en las aproximaciones de las probabilidades posteriores al modelo verdadero y los modelos aproximados.

## TEORÍA DE INFORMACIÓN.

A continuación se presenta la investigación de Pericchi (1984), trabajo que considera el planteamiento y desarrollo del procedimiento bayesiano utilizado por Zellner (1971) y Atkinson (1978). Renombra este procedimiento por el título de Procedimiento Bayesiano Estándar y lo maneja en términos de cantidad de información esperada acerca de un vector aleatorio.

Debido al análisis de diversas situaciones del procedimiento estándar en que se concluyen los resultados arbitrarios, Pericchi encuentra explicaciones y elabora una modificación en su estudio por medio de la cantidad relativa a la ganancia esperada en información.

Con las suposiciones habituales de este procedimiento estándar, se establecen; la existencia de un único modelo verdadero en los modelos alternativos,  $M_j : Y = X_{nj} \theta_j + \epsilon$  ( $j = \overline{1, J}$ ), que los errores en los modelos son no correlacionados.

Para el caso de contemplar modelos con  $\sigma^2$  conocida, se asume tener una densidad inicial Normal Multivariada para el parámetro,  $p(\theta_j/h)$  tiene media  $\theta_{0j}$  y precisión  $V_{0j}^{-1}h$  donde  $h = \sigma^2$  (es decir,  $p(\theta_j/h) \sim N_m(\theta_{0j}; V_{0j}^{-1}h)$ ,  $k_j$  es la dimensión de  $\theta_j$  y la probabilidad inicial del  $j$ -ésimo modelo, es  $\frac{1}{m}$  ( $P_r(M_j) = w_j = 1/m$ ).

Pericchi plantea el procedimiento estándar ya presentado anteriormente:

$$J_{mi} = \left( \frac{R_{ni} + R_{nm} + Q_{ni} + Q_{nm}}{2\delta^2} \right) + \frac{1}{2} \log \left( \frac{|V_{om}^{-1}| |V_{oi}^{-1} + V_{ii}^{-1}|}{|V_{om}^{-1} + V_{im}^{-1}| |V_{oi}^{-1}|} \right), \quad (2.5)$$

donde

$$\theta_{nj} = (\hat{\theta}_{nj} - \theta_{oj})^t (V_{oj} - V_{ij})^{-1} (\hat{\theta}_{nj} - \theta_{oj}), \quad j = \overline{m, i}$$

$$V_{ij}^{-1} = X_{nj}^t X_{nj}, \quad R_{nj}, \quad Q$$

$\hat{\theta}_{nj}$  son las sumas de cuadrados de los residuales estandarizados de

$\hat{\theta}_j$  es el estimador de mínimos cuadrados de

El segundo término de (2.5) es igual a la diferencia entre la ganancia esperada en información acerca de los parámetros entre los modelos, esto es

$$\frac{1}{2} \log \left( \frac{|V_{om}^{-1}| |V_{oi}^{-1} + V_{ii}^{-1}|}{|V_{om}^{-1} + V_{im}^{-1}| |V_{oi}^{-1}|} \right) = I^{\theta_i} \{ \epsilon, p(\theta_i) \} - I^{\theta_m} \{ \epsilon, p(\theta_m) \}, \quad (2.6)$$

donde análogamente para:

$$I^{\theta_i} \{ \epsilon, p(\theta_i) \} = H \{ p(\theta_i) \} - \int p(y) H \{ p(\theta_i; y) \} dy$$

$$H \{ p(\theta_i) \} = - \int p(\theta_i) \log \{ p(\theta_i) \} d\theta_i$$

$$p(y) = \int p(y|\theta_i) p(\theta_i) d\theta_i$$

Considerando a  $\sigma^2$  desconocida, se hace la suposición que la conjugada inicial bajo  $M_j$ , tiene densidad Normal Gamma y se escribe como:

$$p(\theta_j, h_j) = p(\theta_j | h_j) p(h_j) \quad (2.7)$$

Esta última expresión, Perichí la plantea en términos de información esperada, como sigue:

$$I^{\theta, h} \{ \epsilon, p(\theta, h) \} = I^{\theta} \{ \epsilon, p(\theta | h_j) \} + I^h \{ \epsilon, p(\theta | h) \}. \quad (2.8)$$

Y además, la densidad de  $h_j$  ( $h_j = \sigma_j^{-2}$ ) es una Gamma con parámetros  $d_j$  y  $\delta_j$ . En base a esto, se obtiene un procedimiento estándar para  $\sigma^2$  desconocida.

$$S_{m_i} = \log \left\{ \frac{\delta_{m_i}^{d_{m_i}} \delta_{n_i}^{d_{n_i}} \Gamma(d_{m_i}) \Gamma(d_{n_i})}{\delta_i^{d_i} \delta_{nm}^{d_{nm}} \Gamma(d_{nm}) \Gamma(d_{n_i})} \right\} + \frac{1}{2} \log \left\{ \frac{|V_{0i}^{-1} + V_i^{-1}| |V_{0m}^{-1}|}{|V_{0i}^{-1}| |V_{0m}^{-1} + V_{im}^{-1}|} \right\}, \quad (2.9)$$

donde

$$\delta_{n_j} = \frac{1}{2} (2\delta_j + R_{n_j} + O_{n_j}), \quad \bar{j} = \bar{n}_j \quad y$$

$$d_{n_j} = d_j + \frac{1}{2}.$$

Perichí dice que bajo la suposición (2.7) y (2.8) se obtiene

$$I^{\theta, h} \{ \theta, p(\theta, h) \} = \log \left\{ \frac{|V_0^{-1} + V_i^{-1}|^{1/2} \Gamma(d)}{|V_0^{-1}|^{1/2} \Gamma(d_n)} \right\} + \alpha_n \{ \Psi(d_n) - 1 \} - \alpha \{ \Psi(d) - 1 \}, \quad (2.10)$$

donde

$\Psi(x) = \frac{d}{dx} \log \Gamma(x)$  es la función digamma

y usando [2.10], se expresa nuevamente [2.9]

$$S_{mi} = n \log \left( \frac{\partial \pi_i}{\partial \theta_m} \right) + I^{\theta_m, h} \{ \xi, p(\theta_i, h) \} - I^{\theta_m, h} \{ \xi, p(\theta_m, h) \} \quad (2.11)$$

El autor concluye que estos resultados ocasionan por su definición diversos problemas, como por ejemplo, la expresión (2.6) crea engañosas regiones en el espacio paramétrico<sup>(1)</sup> y [2.11] plantea resultados inconsistentes.

Perichi propone una solución a este problema, que consiste en asignar pesos iniciales a los modelos sin discriminar a ningún modelo ( $P_r(w_j) = 1/3$ ). Define al vector de parámetros desconocidos como  $\underline{\eta}$  y a  $M$ , una variable aleatoria que toma  $J$  modelos ajenos y simples (exclusivos). El objetivo principal es maximizar la pérdida de información acerca de  $M$  y maximizar lo que se espera ser aprendido acerca de los parámetros de los modelos.

Estableciendo simbólicamente lo que se espera ser aprendido con respecto a  $\underline{\eta}$  sobre la distribución posterior, se establece para  $\underline{w} = (w_1, \dots, w_j)$

(1) Esto se refiere a que el procedimiento estándar infla la probabilidad posterior de los modelos, que tienen un incremento de información esperada pequeño acerca de sus parámetros.

$$R\{\varepsilon, p(\eta, M)\} = \sum_j w_j [I^{\eta_j}\{\varepsilon, p(\eta_j)\} - \log w_j]. \quad (2.12)$$

Para la maximización de (2.12) se toma  $\tilde{w} = (\tilde{w}_1, \dots, \tilde{w}_j)$ , que es,

$$\max_w R\{\varepsilon, p(\eta, M)\}. \quad (2.13)$$

Esta solución, se expresa como:

$$\log \tilde{w}_i - \log \tilde{w}_k = I^{\eta_i}\{\varepsilon, p(\eta_i)\} - I^{\eta_k}\{\varepsilon, p(\eta_k)\}, \quad (2.14)$$

para modelos ajenos y simples

Algunas observaciones sobre este criterio son: si  $I^{\eta_i}\{\varepsilon, p(\eta_i)\} < I^{\eta_k}\{\varepsilon, p(\eta_k)\}$ , en base a (2.14) se penaliza a  $M_i$ , si la ganancia esperada en información acerca de los parámetros es igual para todos los modelos en competencia, el criterio (2.14) coincide con el procedimiento estándar, ya sea (2.5) o (2.11) dependiendo del caso, y más aún, cuando todos los parámetros son exactamente conocidos, la razón de los momios posteriores es igual a la razón de verosimilitudes.

Procedimientos de Broemeling (para el caso de inferir sobre los parámetros).



Cuando se tiene el objetivo de inferir sobre los parámetros, un criterio de selección es el propuesto por Broemeling (1985, pp. 97-99), quien indica seleccionar el modelo con mayor determinante de los coeficientes de regresión en la matriz de precisión de la distribución posterior marginal.

La matriz de precisión de la distribución marginal de  $\theta$  es

$$H = (n-p) X^t X (Y^t Y - Y^t X (X^t X)^{-1} X^t Y)^{-1}, \quad (2.15)$$

donde  $p$  es el número de parámetros en el modelo.

En la siguiente sección se presentarán otros procedimientos ubicados en la elección de modelos con objetivos predictivos.

## 2.2.2 MÉTODOS QUE ESPECÍFICAN LA PREDICCIÓN.

### 2.2.2.1 PROCEDIMIENTOS QUE NO UTILIZAN FUNCIONES DE PÉRDIDA.

El siguiente trabajo que se describe lo desarrollaron conjuntamente Geisser y Eddy (1979), en donde al parecer los resultados tienen un amplio alcance operacional y resulta muy sencillo de aplicarse con la ayuda de las computadoras. Los antecedentes de este método están en los artículos de Stone (1974, 1976) y Geisser (1975), quienes al mismo tiempo y de manera independiente trabajaron sobre la misma idea de Selección de Variables, por lo cual ambos nombran a sus métodos de diferente manera. Stone a su procedimiento lo nombra como validación cruzada (CROSS-VALIDATION) y Geisser lo llama como reuso de la muestra predictiva (Predictive sample reuse, PSR) aunque ambas significan lo mismo.

Como el problema de regresión se tiene que las densidades predictivas de  $\underline{y}$  dado  $X$  son sensibles a los datos, se introduce el procedimiento de validación cruzada (o PSR) para eliminar la observación  $(x_i, y_i)$  para juzgar la potencia predictiva del modelo con la densidad predictiva de  $\underline{y}$  dado  $X$  en el valor  $x_i$ ;  $S_i = S - \{(x_i, y_i)\}$ , con  $S$  igual a  $\{(x_i, y_i) \mid i = \overline{1, n}\} = S$ . Y se concluye con una función que depende del modelo  $\alpha$  ( $\alpha = \overline{1, m}$ ), que se plantea así:

$$A(\alpha) = \sum_{i=1}^n \log f(y_i | x_i, \alpha, S_i) \quad (2.16)$$

En el artículo de Stone 1976, se muestra que el criterio de validación cruzada es asintóticamente equivalente bajo ciertas condiciones al criterio de Akaike (esto es, un criterio para seleccionar modelos).

El criterio de Akaike estriba en seleccionar al modelo  $M$  que maximice a

$$L(M_\alpha, \hat{\theta}_{M_\alpha}) - P_{M_\alpha}, \quad \alpha = \overline{1, m}. \quad (2.17)$$

donde  $L(M_\alpha, \hat{\theta}_{M_\alpha})$  es el logaritmo de la función de verosimilitud,  $\hat{\theta}_{M_\alpha}$  es el estimador máximo verosímil del parámetro  $\theta_{M_\alpha}$  en el modelo  $M_\alpha$  y  $P_{M_\alpha}$ , el número de parámetros. Por otro lado, se tiene que el criterio de Akaike es equivalente si los modelos tienen varianzas conocida, a otro procedimiento de selección en modelos de regresión lineal múltiple, que consiste en minimizar la C'p de Mallows (Draper-Smith, 1981, pp.299), esta es:

$$C_p = \frac{SCR_{M_\alpha}}{\sigma^2} - (n - 2P_{M_\alpha}) \quad \alpha = \overline{1, m} \quad (2.18)$$

donde  $SCR_{M_\alpha}$  es la suma de cuadrados residuales para el modelo  $M_\alpha$ ,  $n$  el tamaño de la muestra y  $\sigma^2$  la varianza conocida.

Continuando con el trabajo de Geisser y Eddy (1979), estos también omiten la posible influencia de la  $i$ -ésima observación de la muestra cuando se calcula la densidad predictiva de  $X_i$ . Además, se denomina la validación cruzada mediante el nombre de predictive sample reuse (PSR), y se deriva un criterio que estriba en maximizar el producto de predictivas condicionales.

Se define el vector de datos  $X_{(i)} = (x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$  con la  $i$ -ésima observación eliminada, con base en esto, la densidad predic

tiva de  $x_i$  cuando  $M_A$  es el modelo verdadero, se representa como  $f_i(x_i/x_{(i)}, S_i, M_A)$ . En base a esto, se da el criterio de selección que consiste en elegir aquella  $M_A$  que minimice  $L_A$ , eligiendo entonces al modelo  $M_{A^*}$ , donde

$$L_A = \prod_{i=1}^n f_i(x_i/x_{(i)}, S_i, M_A) \quad , \quad d = \bar{h}m. \quad (2.19)$$

Se presentan dos maneras distintas para obtener la densidad predictiva de observaciones futuras, cuando  $M_A$  es verdadero, que se sustentan en: 1) como  $f(x/\theta_A, S_i, M_A)$  es la densidad de las observaciones y  $\hat{\theta}_A$  es el estimador máximo verosímil, calcular  $f(x/\hat{\theta}_A, S_i, M_A)$ , esta densidad es llamada "densidad predictiva cuasi-verosímil" y 2) calcular la densidad predictiva bayesiana obtenida de una distribución inicial vaga.

El siguiente resultado, es un caso particular de este procedimiento, pues sólo es válido para el anidamiento de modelos, esto es, si  $M_A$  está anidado en  $M_{A'}$  ( $M_A \in M_{A'}$ ), se tiene que la densidad predictiva cuasi-verosímil bajo  $M_A$  es  $f_i(x_i/S_i, \hat{\theta}_{A(i)}, M_A)$ , donde  $\hat{\theta}_{A(i)}$  es el estimador máximo verosímil de  $\theta_A$  omitiendo a  $x_i$ , entonces si además  $M_A$  es verdadero, el cociente  $L_A/L_{A'}$  converge a

$$\frac{-p(M_A) + o(M_A) + \log \lambda}{-p(M_{A'}) + o(M_{A'}) + \log \lambda} \quad , \quad \text{si } n \rightarrow \infty. \quad (2.20)$$

donde  $p(M_A)$  y  $p(M_{A'})$  es el número de parámetros desconocidos en  $\theta_A$  y  $\theta_{A'}$ , respectivamente

$$\lambda = \frac{\prod_{i=1}^n f_i(x_i/S_i, \hat{\theta}_A, M_A)}{\prod_{i=1}^n f_i(x_i/S_i, \hat{\theta}_{A'}, M_A)}$$

Para este caso, los autores concluyen diciendo que (2.20) es esencialmente el criterio de Akaike. Si  $\lambda^1$  se elige al modelo  $M_{\lambda^1}$ , y en caso contrario ( $\lambda^1$ ) se elige al modelo  $M_{\lambda^2}$ .

Como última contribución, los autores trabajaron otras situaciones, donde no se cuenta con una función de verosimilitud para al menos uno o más de los modelos involucrados, por lo que se introduce una medida de la discrepancia puntual entre  $x_i$  y  $\hat{x}_i$ , con  $i = \overline{1, n}$ , dada por  $D_n = \sum_{i=1}^n d(x_i, \hat{x}_i)$ , donde  $d(x_i, \hat{x}_i)$  es alguna medida de la distancia entre  $x_i$  y  $\hat{x}_i$ . Este resultado, se generaliza para los modelos de regresión múltiple, y se concluye eligiendo al modelo que tenga una discrepancia menor.

Broemeling (1985, pp.97-99) también aporta esta sección de métodos predictivos un procedimiento que involucra a la distribución predictiva de las observaciones futuras ( $P_k(y_0/S)$ ), donde su limitación depende en conocer previamente el valor a predecir,  $x_0$ ;  $S$  es la muestra y  $y_0$  la observación futura,  $k$ , el  $k$ -ésimo modelo y  $P_k(y_0/S)$  se distribuye como una T-Student, para los  $m$  modelos en general. El criterio consiste en elegir el modelo que haga máxima la distribución predictiva en las observaciones futuras. O sea, elegir  $M_{\lambda^k}$  si

$$P_1(y_0/S) > \max_{m \neq k} P_m(y_0/S), \quad 1 \leq k \leq m. \quad (2.21)$$

Los procedimientos que toman en cuenta explícitamente el problema de predicción en el Análisis de Regresión son los procedimientos que a continuación se presentan.

## 2.2.2.2 TEORÍA DE DECISIONES.

En esta sección se presentan cuatro trabajos desarrollados en una misma línea bayesiana, que radica en plantear el problema de selección de variables con objetivos predictivos en base a un problema de decisión. Estos trabajos son realizados por Lindley (1968), Brooks (1972, 1974) y Goldeisten (1976). Los tres primeros tienen un desarrollo común dentro de la teoría de decisiones y son consecutivos en su realización. En el último artículo se emplea distribuciones iniciales no paramétricas.

El trabajo de Lindley (1968) plantea las bases para tratar el problema de selección de variables por medio de la teoría de decisiones. Se denotará por  $H$ , a todo el conocimiento disponible para tomar la decisión de cuál es el mejor modelo incluyendo a los datos, tanto de las variables independientes como de la variable dependiente.

Enseguida se presentan las suposiciones estructurales, con las que Lindley desarrolla el trabajo. Para algunas de ellas su formulación resulta natural, debido a que van de acuerdo con las suposiciones usuales del Análisis de Regresión.

Suposición 1:  $\theta$  y  $X$  son variables aleatorias en  $\mathbb{R}^r$  (espacio euclídeo de  $r$  dimensiones) las cuales son independientes condicionalmente a  $H$ . Esto es, su densidad conjunta se puede factorizar de la siguiente manera:

$$p(\theta, X | H) = p(\theta | H) p(X | H)$$

(2.22)

Suposición 2:  $\underline{y}$  es una variable aleatoria en  $\mathbb{R}^n$  con densidad  $p(\underline{y}/\theta, x, H)$ , la cual no depende de  $H$ . Es decir  $p(\underline{y}/\theta, x, H) = p(\underline{y}/\theta, x)$  y su densidad tiene la propiedad de

$$\begin{aligned} E(\underline{y}/\theta, x) &= \theta^t x \\ \text{Var}(\underline{y}/\theta, x) &= \sigma^2 \end{aligned} \quad (2.23)$$

Las siguientes suposiciones presentan restricciones fuertes en su planteamiento.

Suposición 3:  $\sigma^2$  es conocida. Se denotará por  $I$  a un subconjunto (puede ser vacío) de los enteros  $1, 2, \dots, r$  tal que contenga  $s$  elementos y a su respectivo complemento, que tiene  $r-s$  elementos, se le nombrará mediante  $J$ .

Sea  $X$  tal que  $X = (X_1, \dots, X_r)$  y se define a  $X_I$  de tal forma que sea el vector cuya componente  $X_i$  es tal que  $i \in I$ , análogamente se define  $X_J$ . O sea, mediante  $X$  se denota a las variables explicativas y con  $I$  se determina a un subconjunto  $X_I$  de éstas.

Suposición 4: El espacio de decisión consiste de elementos  $d = (I, f(\cdot))$ ; donde  $I$  ya se describió anteriormente y representa los subconjuntos de variables explicativas que se van a considerar y por  $f(\cdot)$ , a una función de  $\mathbb{R}^s$  a  $\mathbb{R}$  que es llamado el "predicador de  $y$ " porque evalúa la predicción  $f(x_I)$ , en todos los subconjuntos de variables explicativas.

Suposición 5: La función de pérdida es:

$$\{y - f(x_I)\}^2 + c_I \quad \text{con} \quad c_I \geq 0, \quad (2.24)$$

donde  $y$  es el valor por predecir de la variable independiente<sup>(1)</sup>.

Para comprender esta última suposición se presenta la anatomía o estructura cualitativa de un problema de selección de variables, como un arreglo cronológico en un árbol de decisión



Este árbol cuenta en su diagrama con nodos cuadrados y circulares. Los primeros nodos denotan situaciones de elección o decisión, y los segundos simbolizan situaciones que están influidas bajo un proceso aleatorio. El árbol se lee de izquierda a derecha. La figura anterior representa un problema de selección de variables con objetivos predictivos, el primer nodo representa el conocimiento particular del marco de referencia en que se desenvuelve el problema, así como la decisión de cuales serán las variables explicativas que se van a observar. Ya teniendo claro cuales son las variables explicativas involucradas (en I) para el problema, estas se observarán en condiciones generales y se obtendrán los valores,  $x_I$ . Para estos valores se calcularían las predicciones,  $f(x_I)$ , y se llega al problema de decidir cual es la mejor predicción de los subconjuntos de variables explicativas que predicen a  $y$ . Finalmente considerando al verdadero valor de  $y$ , se evalúa la función de pérdida descri

(1) Entendiéndose por esto, que  $y$  es un valor que va a ocurrir (6 ya paso).



ta en (2.24) y con la función de pérdida se describe la exactitud de la predicción realizada.

Técnicamente esta situación se resuelve de derecha a izquierda para analizar su desarrollo y planteamiento, esto es, si se minimizó la pérdida esperada en (2.24) quiere decir, que la diferencia entre  $y$  y  $f(x_2)$  fue muy pequeña de donde,  $f(x_2)$  es un buen predictor que está siendo obtenido por un buen subconjunto de variables explicativas y a las que además, se les tomó en cuenta su costo de obtención.

Las suposiciones 4 y 5 sustentan el desarrollo del trabajo de Lindley en cuanto a la teoría de decisiones. A continuación se presentarán otras tres suposiciones utilizadas en el desarrollo de Lindley.

Suposición 6: La distribución  $p(X/H)$  (que cumple la suposición 1) tiene la siguiente característica

$$E(X_2/X_1) = \beta_{JX} X_1 \quad (2.25)$$

con

$\beta_{JX}$  matriz de  $(r-s) \times s$

Suposición 7: Ahora  $H$  contiene  $n$  observaciones estocásticamente independientes  $y_1, \dots, y_n$  con valores en  $x_{ij}$ ; ( $i = \overline{1, n}; j = \overline{1, n}$ ) de las variables independientes con densidad  $p(y_i/\theta, x_i)$  dadas por la suposi-

ción 2 y se tienen las siguientes relaciones

$$E(\underline{y}/\underline{x}, \underline{\theta}) = \underline{x} \underline{\theta} \quad \text{y} \quad V(\underline{y}/\underline{x}, \underline{\theta}) = \sigma^2 \underline{I}_n \quad (2.26)$$

Suposición 8: Las  $\underline{X}_i$  (de la suposición 7) son variables aleatorias idénticamente distribuidas según una normal multivariada y además, se supone independencia de la variable aleatoria  $\underline{X}$ , vector en el cual se realiza la predicción con los valores  $\underline{X}_i$ . Los parámetros de la distribución de  $\underline{X}$  son independientes de  $\underline{\theta}$ .

El siguiente supuesto establece el manejo de distribuciones iniciales de referencia.

Suposición 9: Se denotará por  $\mathcal{H}_0$  a la información inicial que plantea el interesado. Dado este conocimiento,  $\mathcal{H}_0$ , la distribución de  $\underline{\theta}$  es uniforme sobre  $\mathcal{R}^n$ , la distribución de las medias, de la distribución normal de  $\underline{X}$  es similarmente uniforme sobre  $\mathcal{R}^r$ , la distribución de la matriz de dispersión de la distribución normal de  $\underline{X}$  es inversamente proporcional a su determinante. Y además se pide que estas tres distribuciones sean independientes.

Bajo las anteriores suposiciones, Lindley establece un teorema en donde se obtiene la solución óptima del problema de predicción, que consiste en elegir el subconjunto de variables independientes  $I$  que satisfaga

$$\min_I \left\{ \frac{R(\underline{J}:\underline{I})}{n} + C_x \right\} + \sigma^2 \left( 1 + \frac{r}{n} \right)$$

y se predice a  $y$  mediante  $E(\underline{\theta})^t E(x/x_I)$

donde

$$R(J:I) = E(\underline{\theta}_j)^t V(x_j) E(\underline{\theta}_j),$$

$R(J:I)$  es la reducción en suma de cuadrados debida a  $\underline{\theta}_j$  fijando  $\underline{\theta}_I$ ,  $V(x_j)$  es la partición de  $V(Y/H) \sim \frac{w^t w}{n}$ <sup>(1)</sup> y está, es la matriz de varianzas-covarianzas.

Eliminando los términos constantes en (2.27), se encuentra que la solución óptima para el problema de predicción radica en seleccionar a la  $I$  que satisfaga

$$\min_I \left\{ \frac{R(J:I)}{n} + C_I \right\}. \quad (2.28)$$

Cuando este resultado se compara con otros criterios o técnicas, por ejemplo con la estadística  $C_p$  de Mallows<sup>(2)</sup>, se encuentra una relación muy cercana entre ambas técnicas bajo una transformación lineal si se toman costos iguales. O sea, si por  $R_I$  se denota a la suma de cuadrados residual fijando las variables  $x_i$ ,  $i \in I$ , se define como

$$C_p = \frac{R_I}{\sigma^2} - (n - 2p), \quad (2.29)$$

(1)  $V(Y/H) = \frac{w^t w}{n(n-2)}$  donde  $w$  es una matriz cuyo elemento típico es  $x_{ij} - x_{.j}$ . La matriz aproximadamente es igual a la matriz de dispersión muestral.

(2) Esta estadística ya fue presentada en el artículo de Geisser y Eddy (1979).

donde  $\sigma^2$  es la varianza estimada,  $R_I = R(J:I) + R$ ,  $R$  es la estimación residual de todas las variables.

Por lo que basándose en la suposición (que es muy restrictiva), en donde cada una de las variables tienen el mismo costo ( $C=2\sigma^2/n$ ) se obtiene.

$$C_p \frac{\sigma^2}{n} = \left\{ \frac{R(J:I)}{n} + c_I \right\} + \frac{R}{n} - \sigma^2. \quad (2.30)$$

Como ejemplos, se presenta el procedimiento aplicado a los datos de Hald que tienen cuatro variables explicativas. Se les trabajó con la suma de cuadrados residuales y suponiendo aditividad en sus costos, se elige en cada caso el mínimo de ellas. Es decir utilizando la expresión más simple de la solución del problema de predicción (2.28). La decisión óptima es consecuencia de elegir el modelo que minimice la siguiente expresión:

$$\begin{aligned} \min_I \left\{ \frac{R(J:I)}{n} + c_I \right\} &= n \min_I \{ R_I - R + n c_I \} \\ &= \min_I \{ R_I + n c_I \}. \end{aligned}$$

TABLA No.1 ANALISIS DE PREDICCIÓN EN LOS DATOS DE HALO

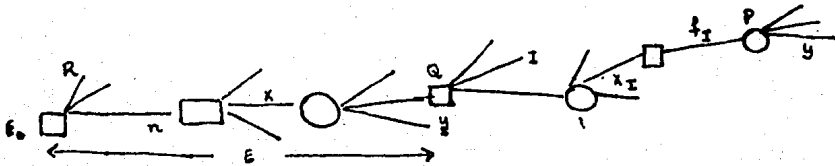
Variables en I	$R_T$	$R_{I+nC_T}$	$R_{I+nC_T}$	$R_{I+nC_T}$
ninguno	2716	2716	2716	2716
1	1266	1766	1366	2266
2	906	1406	1006	1006
3	1939	2039	2439	2039
4	884	984	1384	2884
12	58	1058	258	1158
13	1227	1827	1827	2327
14	75	675	675	3075
23	415	1015	1015	615
24	869	1469	1469	2969
34	176	376	1176	2276
123	48 <sup>+</sup>	1148	748	1248
124	48 <sup>-</sup>	1148	748	3148
134	51	751	1151	3151
234	74	774	1174	2274
<b>todos</b>	<b>48</b>	<b>1248</b>	<b>1248</b>	<b>3248</b>
<u>costos</u>	$nC_1$	500	100	1000
	$nC_2$	500	100	100
	$nC_3$	100	500	100
	$nC_4$	100	500	2000

Considerando ahora los trabajos de Brooks (1972 y 1974) se encuentra que se hizo parte del procedimiento de Lindley, planteando ahora el problema de encontrar un diseño óptimo de regresión, para encontrar las mejores variables para poder predecir. Esto es, realizar un diseño experimental que además de lo que encontró Lindley, se tome en cuenta que variables se van a observar, en que rango de variabilidad

se van a medir, en que momento se van a observar, costos de observación, costos de obtención, costos de predicción y otros más.

Un aspecto importante del trabajo de Brooks, se debe a la ampliación del trabajo de Lindley para cuando  $\sigma^2$  es desconocida, que prácticamente consiste en estimar a  $\sigma^2$ .

El arreglo cronológico del árbol de decisión, se plantea así:



En general, el problema de encontrar el diseño óptimo de regresión, es un planteamiento más completo que el discutido por Lindley, resultando que las expresiones involucradas son más complicadas. Sin embargo, sus expresiones finales involucran el criterio de decisión de Lindley.

La extensión de  $\sigma^2$  desconocida, Brooks la solución en base a la suposición 1, pues ahí se asumió que el vector de parámetros  $(\theta, h = 1/\sigma^2)$  y  $X$  son variables aleatorias, dado el conocimiento inicial,  $E_0$ . Mediante lo anterior dice que lo importante es reemplazar a  $\sigma^2$  por  $E(\sigma^2/E_0)$ , donde  $E(\sigma^2/E_0)$  representa el valor esperado de la distribución inicial de  $\sigma^2$ .

El criterio de decisión consiste en elegir el modelo que haga mínima la siguiente expresión:

$$\min_{\mathbf{I}} \left\{ \omega + \text{tr} \left[ \omega \mathbf{M} (\mathbf{z}' \mathbf{z})^{-1} \right] + c_{\mathbf{I}} + c \right\} \quad (2.31)$$

donde

$$\omega = E(\sigma^2/\epsilon_0), \mathbf{M} = E(\mathbf{z}\mathbf{z}'/\epsilon_0), \mathbf{z}' = (\mathbf{1}, \mathbf{X}')$$

$C$  denota el costo del experimento por unidad de error en la predicción  $C(n, \mathbf{X})$  y  $c_{\mathbf{I}}$ , denota el costo de las observaciones de las variables  $x_{\mathbf{I}}$  por costo de unidad de error en la predicción.

Goldstein (1976) plantea su expresión sobre el problema de regresión bayesiano usando una distribución no paramétrica para el parámetro desconocido  $\theta$ . El manejo y planteamiento de las distribuciones iniciales no paramétricas se encuentran en un trabajo de 1975 de Goldstein.

Tomando en cuenta todas las consideraciones dadas en el problema de regresión (Capítulo 1), el autor señala la no justificación que se hace siempre al pedir entre las variables explicativas y la variable dependiente una relación lineal (esto es,  $E(\mathbf{y}/\mathbf{X}=\mathbf{x}) = \mathbf{X}\beta, \mathbf{y}_1$ ), por lo que propone realizar una correcta selección de  $\theta$ , considerando un buen predictor de  $\mathbf{y}$ , de esta manera para el caso de tener a  $\sigma^2$  conocida se establece un criterio razonable para obtener un estimador bayesiano de  $\theta$ , que consiste en seleccionar el valor que minimice el riesgo  $\mathcal{R}(\beta)$ , donde

$$\begin{aligned} R(\beta) &= E(y - \sum \beta_i x_i)^2 \\ R(\beta) &= E((y - \sum \beta_i x_i)^2 / y, x, \beta). \end{aligned} \quad (2.32)$$

Este criterio da el mejor predictor lineal de  $Y$  dado  $X$ . Los valores óptimos,  $\beta^*$ , que minimizan  $R(\beta)$  y  $R(\beta^*)$ , están dados por:

$$\beta^* = D^{-1}b \quad \text{y} \quad R(\beta^*) = \frac{\begin{vmatrix} D & b \\ b' & c \end{vmatrix}}{|D|} \quad (2.33)$$

con

$$\begin{aligned} D &= (d_{ij}) \quad d_{ij} = E(x_i x_j / y, x, \beta), \quad i, j = \overline{1, p}, \\ b &= (b_1, \dots, b_p)' \quad b_i = E(y x_i / y, x, \beta) \quad \text{y} \quad c = E(y^2). \end{aligned}$$

Goldstein prueba que para  $n$  grande, el estimador de mínimos cuadrados (capítulo 1) es una buena aproximación para  $\beta^*$  sobre un rango amplio de distribuciones iniciales del parámetro desconocido y se dice que este resultado es de gran ayuda cuando el cálculo de (2.33) sea muy difícil.

Considerando el problema de seleccionar el óptimo subconjunto de variables de regresión se propone realizar la comparación del riesgo expresado en  $R(\beta^*)$ . Este procedimiento no especifica el tamaño a partir de cual, el riesgo es inaceptable, por lo que se espera que el interesado o estadístico lo determine en base a sus necesidades.

El planteamiento para saber si un subconjunto de variables puede ser eliminado es el siguiente: Se define  $S_k$  un conjunto de elementos  $\{i_1, i_2, \dots, i_k\}$  y la matriz  $D(S_k)$  de tamaño  $(p-k) \times (p-k)$  que es el resultado de la reducción de renglones y columnas de la matriz  $D$  y  $b(S_k)$ .



vector de  $(p-k)$  que proviene de la reducción de  $i_1, \dots, i_k$  entradas en  $D$ , entonces

$$R_{S_k}(\beta^*) = \frac{\begin{vmatrix} D(S_k) & b(S_k) \\ b^*(S_k) & c \end{vmatrix}}{|D(S_k)|}, \quad (2.34)$$

$$\Delta = R_{S_k}(\beta^*) - R(\beta^*). \quad (2.35)$$

Por  $\Delta$  se denota al riesgo incurrido si omitimos las variables  $(x_{i_1}, \dots, x_{i_k})$  que significa el incremento en el valor mínimo de (2.32), si se omiten las variables y también se pueden eliminar subconjuntos de variables para las cuales  $\Delta$  no sea mayor que un valor determinado anticipadamente. Este valor fijo establecerá un compromiso entre las demandas de un buen modelo para predecir y la simplicidad para predecir.



**CONCLUSIONES**

## CONCLUSIONES

A continuación se presentan las conclusiones que se desprenden del anterior capítulo referentes a los procedimientos de selección de variables, las cuales son: procedimientos de probabilidades posteriores, procedimiento de Broemeling, procedimiento de Geisser y Eddy, procedimiento de Lindley y procedimiento de Goldstein.

Se comenzará comentando una dificultad global referente a la funcionalidad de los procedimientos de selección de variables, ya sea tomando el punto de vista Clásico o Bayesiano. Esta dificultad se presenta en todos los procedimientos de selección que se basan en el cálculo de "todas" las regresiones; como el número de posibles modelos a calcular depende del número de variables explicativas (a saber  $p$ ), el total de posibles modelos a calcular dependerá de la cantidad de variables explicativas con las que se cuenten, que en general serán  $2^p - 1$  modelos distintos. En la práctica puede ocurrir que el número de variables explicativas sea grande, por lo que no es recomendable aplicar de inmediato alguno de los procedimientos antes mencionados de selección, el adecuado a juicio del interesado o el estadístico, sino que se sugiere realizar una "selección" de variables previa, por medio de un criterio extramatemático o un criterio de selección (por ejemplo *backware* o *stepwise*) que al no presentar este problema, permite eliminar algunos de los modelos menos importantes (dependiendo del caso a analizar) y permite obtener una idea sobre el conjunto de variables "óptimo". Con esta información en mente se aplicará alguno de los procedimientos de selección estudiados en este trabajo. De no seguir esta recomendación, se tendrá que usar una cantidad impresionante de tiempo de máquina, pues los resultados requerirán de tal vez años de uso de computadora.

Tomando como referencia los resultados que Atkinson obtuvo, se vé que este criterio no es posible de aplicar tal cual a los modelos de regresión, ya que se tienen graves consecuencias al asignar las probabilidades a la mayoría de los modelos de regresión, porque éstos en general son modelos anidados y al darle "un peso" a un subconjunto de variables explicativas, se les tendrá que asignar el mismo o mayor peso dependiendo del caso a todos los modelos que involucren ese conjunto de variables, ocasionando así, que la suma de todos los modelos pudiera ser mayor que uno, lo cual es una contradicción a los supuestos de la teoría de probabilidades. Otro resultado desfavorable, se tiene en la asignación de las probabilidades de los modelos pues modelos formados por ciertos subconjuntos de variables explicativas se les asignará un peso, que va hacer menor que el peso de cualquier otro modelo formado por más de esas variables explicativas y esta probabilidad va hacer menor que la probabilidad asignada al modelo completo. Esto ocasionaría que el modelo completo tendrá mayor oportunidad de ser elegido, que otros modelos. Por lo cual este criterio siempre eligirá al modelo completo y de aquí, que sea un mal procedimiento de selección.

Otra dificultad que se tiene en este procedimiento de probabilidades posteriores, es la restricción de aplicarlo solo a modelos con el mismo número de parámetros, por lo que en caso de utilizarse en los modelos de regresión se tiene que hacer una clasificación en niveles, de tal forma que cada uno de ellos se agrupe a los modelos con igual número de parámetros. Entonces se presenta otro problema al tratar de asignar la probabilidad de los modelos, pues en cada uno de ellos se va asignar la probabilidad condicionada al nivel y a las variables explicativas que se tengan, Esto hace más difícil la aplicación de este modelo.

Se concluye lo siguiente, este procedimiento no es recomendable.

Ejemplificando lo anterior, si el número de variables explicativas es 10 se tendrán que calcular  $2^{10} - 1 = 1023$  modelos diferentes; y si por cada modelo, la máquina tarda un promedio de 1.5 minutos, se tendrán que esperar 4534.5 minutos, que corresponden a 75.6 horas y en días serían, 3,15 días. Este resultado muestra que si el número de variables es mayor que 10 y considerando que el número de modelos diferentes crece de manera exponencial, entonces se necesitarán más de 3 días de "uso exclusivo" de la computadora y se concluye, que la solución dada por el procedimiento es casi imposible de obtener por requerir mantener tantos recursos ocupados en un solo problema.

A continuación se comentará cada uno de los diferentes procedimientos de selección estudiados.

#### PROBABILIDADES POSTERIORES;

Es importante notar que este procedimiento viola los axiomas de coherencia, por lo cual, la decisión de elegir un modelo determinado no es necesariamente buena.

El procedimiento de probabilidades posteriores inicialmente tiene problemas en lo referente a la asignación de la probabilidad a los modelos, pues la interpretación de la probabilidad no se indica explícitamente en los artículos referidos (Zellner y Atkinson), aunque en otro artículo posterior a estos (Geisser y Eddy, 1979) se indica que se refieren a la probabilidad asignada de que el modelo sea el modelo verdadero. También el problema de asignar probabilidades se agudiza para el caso de trabajar con muchas variables, pues cada uno de los modelos deberá tener asignada una probabilidad y la suma de todas debe valer uno.

para selección de variables en regresión. Quizá resulte válido para modelos que no son de regresión o fuera de un contexto de la selección de variables y por individuos que no toman en cuenta los axiomas de coherencia.

#### PROCEDIMIENTO DE BROEMELING.

Este procedimiento está diseñado para inferir sobre los parámetros desconocidos. Es un método basado en el cálculo de la matriz de precisión de la distribución posterior marginal de cada uno de los modelos. Además, es el único método encontrado cuyo objetivo explícitamente es inferir sobre los parámetros y es además, sencillo de aplicarse.

A continuación se presentan los procedimientos que tienen objetivos predictivos, estos son: los métodos de Broemeling, Geisser y Eddy, Goldeistein y Lindley. En los primeros tres métodos, no importa que los parámetros de regresión sean o no desconocidos. El método de Lindley, sólo se plantea en el caso de tener varianza conocida, sin embargo se puede aplicar considerando la estimación de la varianza. Así lo aplico Lindley en los ejemplos que presenta y así, lo trabajo posteriormente Brooks.

#### PROCEDIMIENTO DE BROEMELING,

Este método se aplica si se desea predecir un valor de la variable dependiente, cuando el vector de las variables explicativas toma un valor en particular. Esta situación por lo general no se presenta en los problemas de regresión y de aquí, que es un factor limitante para la aplicabilidad de este procedimiento.

El procedimiento de Broemeling, consiste en calcular la precisión de la distribución predictiva en la observación futura, para cada uno de los modelos que se tengan y elegir al modelo que maximice la predicción futura o elegir al modelo que minimice la varianza.

Para el caso de tener interés en predecir el valor de la variable dependiente en varios vectores de las variables explicativas, es necesario efectuar el mismo procedimiento todas las veces que sea necesario porque el método calcula la precisión de la distribución predictiva en base en cada una de las observaciones futuras. De donde puede concluirse que si se tiene esta situación se recomienda no usar este procedimiento.

#### PROCEDIMIENTO DE GEISSER Y EDDY.

Este es un procedimiento basado en el criterio de validación cruzada, que consiste en eliminar la posible influencia de la  $i$ -ésima observación de la muestra, al calcular la distribución predictiva tantas veces como observaciones se tengan. Estas densidades a su vez se multiplican entre sí y el cálculo se repite para cada uno de los  $2^k$  modelos de regresión que se van a comparar. De aquí, que sea un método que requiere de un mayor manejo de cómputo que los procedimientos presentados aquí. Para algunos problemas de selección de variables en regresión en que  $p$  es "relativamente pequeño", se tiene un proceso de cómputo relativamente sencillo. Más sin embargo, existen otros problemas en donde es imposible esperar la respuesta del procedimiento, ya que el tiempo requerido es muy grande.

El procedimiento de Geisser y Eddy, es un método que tiene ventajas y desventajas, las cuales son: la principal desventaja es ser un método caro porque requiere un fuerte manejo de cómputo aunque  $p$  sea

"relativamente pequeño" y las ventajas: 1) no se necesita un valor específico para la predicción y 2) la utilización del criterio de validación cruzada permite poner a prueba cada distribución predictiva con base a toda la información y esto, trae como consecuencia elegir la mejor distribución predictiva.

Se concluye que es un método recomendable aunque su manejo de cómputo sea grande.

#### PROCEDIMIENTO DE LINDLEY.

Es el único método encontrado que toma en cuenta explícitamente los costos que se involucran en el proceso de selección de variables de regresión. Este es un criterio planteado con base en la Teoría de Decisiones.

En el procedimiento de Lindley hay planteamientos que no están muy claros en su manejo, como son: la eliminación de las unidades de los costos para "sumarse" a la suma de cuadrados de los residuales y la supuesta "aditividad" entre los costos. De aquí, que la aplicación de este método no se pueda llevar fácilmente a la práctica.

Del ejemplo presentado por Lindley, se puede ver que es posible usar este procedimiento cuando se suponga que el costo implicado en la obtención de las variables explicativas sea cero y también cuando se conozca la relación entre sus costos, por ejemplo:  $c_1 = 2c_2$ ,  $c_2 = c_3$ ,  $c_4 = c_1$ , sin necesidad de conocer el valor exacto.

Un resultado importante que se desprende del artículo de Lindley



es el que se refiere, a la relación que se guarda de los resultados de Lindley con la Cp de Mallows. Como la Cp es un criterio de selección de variables generalmente utilizada por los estadísticos clásicos, se presenta aquí una posible justificación bayesiana de su utilización. Además, ambos resultados coinciden cuando en el criterio de Lindley se manejan costos ceros.

#### PROCEDIMIENTO DE GOLDESTEIN,

Este es un método que teóricamente podría usarse cuando se desconociera la normalidad de los errores de los modelos de selección.

La dificultad principal de este método radica en que no se menciona la manera de calcular el valor con el cual los modelos se van a comparar. Por lo que este procedimiento no es utilizable en la práctica.

Los procedimientos antes estudiados no necesariamente producen los mismos resultados. Para ilustrar este hecho se aplican los procedimientos utilizando los datos de Hald (Draper y Smith, 1982, pp. 629-630).

A continuación se presenta la tabla No.2 basada en los datos de Hald, que ejemplifican algunos de los procedimientos aquí vistos, como son: el procedimiento de PSR (Geisser y Eddy, 1979) ubicado en la columna 1, la Cp de Mallows (Draper y Smith, 1982, pp. 301) colocada en la columna 2, Discrepancias (Geisser y Eddy, 1979) columna 3, y el procedimiento de Lindley (Lindley, 1968) en la columna 4

TABLA No.2 RESULTADOS DE APLICAR LOS PROCEDIMIENTOS DE GEISSER, EDDY Y LINDLEY (sin costos y con el equivalente a la  $C_p$  de Mollows). Datos de Hald.

VARIABLES EN LA REGRESION	PSR (log $L_0$ )	$ C_p - p $	$D_k$	$R_I$ (sin costos)
Ninguna		442.2		2716
1	-65.86	200.5	130.74	1266
2	-64.50	140.5	92.47	906
3	-68.75	313.2	201.26	1939
4	-63.77	136.7	91.86	884
12	-46.23	0.3	7.22	58
13	-66.28	195.1	170.62	1227
14	-48.13	2.5	9.32	15
23	-59.90	59.4	53.98	415
24	-64.67	135.2	112.45	869
34	-54.00	119.4	22.62	176
123	-45.57	1	6.92	48+
124	-45.29	1	6.57	48-
134	-45.71	0.5	7.27	51
234	-48.33	33	11.30	74
1234	-45.80	0	8.49	48

Si se realiza una selección global de todos los modelos involucrados se observa la siguiente (tabla No. 2): en la columna 1, se elige al modelo con las variables  $X_1, X_2, X_4$ , pues éste tiene un valor máximo de PSR; si se observa la columna 2, se elige al modelo con las variables  $X_1, X_2$ , pues se hace mínima la  $C_p$  de Mallows; utilizando la columna 3, se elige al modelo con las variables  $X_1, X_2, X_4$ , pues se minimiza la discrepancia y por último en la columna 4, con costo cero, en el procedimiento de Lindley se elige al modelo con las variables  $X_1, X_2, X_4$ .

Ahora, viendo el criterio propuesto por Draper y Smith que consiste en calcular todas las regresiones y analizar los modelos agrupados en conjunto de modelos con el mismo número de parámetros. La agrupación resultante se puede analizar de diferentes maneras y en cada una de ellas se obtienen diferentes resultados, a pesar de usar un mismo procedimiento de selección. Esto ocasiona una aparente inconsistencia, que es consecuencia de un manejo inapropiado del criterio de elección y del escaso conocimiento de la estructura del problema que se tiene.

Para ejemplificar lo anterior, si el objetivo es ir construyendo un modelo que vaya introduciendo mayor número de parámetros, por medio del criterio de Lindley, se tiene lo siguiente en la agrupación de los modelos: en el nivel uno (agrupación de modelos con un parámetro) se elige al modelo con la variable explicativa  $X_4$  por tener el valor mínimo en ese nivel, para el siguiente nivel (con dos parámetros) se elige al modelo que tiene a la variable  $X_4$ , estos son: los modelos que involucran a las variables  $X_1, X_4$ ,  $X_2, X_4$  y  $X_3, X_4$ . De estos, el valor mínimo lo obtiene el modelo  $X_1, X_4$ . Siguiendo con este mismo mecanismo en el nivel 3 (con tres parámetros) se elige al modelo  $X_1, X_2, X_4$  por tener un valor mínimo, usando el criterio de Lindley y se finaliza quedándose con él.

Usando este mismo ejemplo, pero de manera diferente como es el construir un modelo que inicie con todas las variables y vaya eliminando a las menos importantes con base al criterio de decisión de Lindley, se tiene lo siguiente: En el nivel 4 (cuatro parámetros) se comienza con el modelo completo, se pasa al siguiente nivel inferior y se encuentra con otros posibles modelos a elegir, en donde el modelo elegido en este nivel es el  $X_1, X_2, X_3$ , pues minimiza el criterio en este nivel. Si de nuevo se pasa al siguiente nivel inferior, se elige al modelo  $X_1, X_2$ . Ahora, al tratar de pasar al nivel 1, el modelo que minimiza el criterio es más de dos veces mayor, que el valor obtenido por el modelo  $X_1, X_2$ , por lo cual se elige a este modelo.

Continuando con el mismo ejemplo, si ahora se considera algún tipo de información extramatemática, como es la dada por las correlaciones, que permite ir decidiendo sobre la entrada y salida de las variables explicativas, se tiene que en el caso de estos datos:  $X_1, X_3$  y  $X_3, X_2$  parecen ser altamente correlacionados, de donde la información de una variable puede ser recogida por la otra, entonces esto permite eliminar: al modelo completo, al modelo  $X_1, X_2, X_3$ , a todos los modelos del nivel tres y a los modelos  $X_1, X_3$  y  $X_3, X_2$ . Por lo cual, se puede concluir en elegir al modelo  $X_1, X_2$ . O quizás, elegir a un modelo de una sola variable, si es necesario ó no interesa que el criterio mínimo en el nivel 2 sea aumentado casi dos veces en el nivel 1 al elegir a la variable  $X_3$ .

Estos resultados muestran la diversidad de conclusiones al manejar diferente el criterio propuesto de Draper y Smith. Cabe aclarar que Broemeling, Geisser y Eddy, y Lindley, sugieren no realizar el criterio anterior solamente, si no dar una elección global que permita conocer la estructura global del problema e identificar las causas que ocasionan la aparente inconsistencia.

**ANEXO**

### 1. LA DISTRIBUCION UNIVARIADA STUDENT.

Una prueba aleatoria se distribuye como una Student, si se puede expresar de la siguiente manera

$$p(x/\theta, h, v) = \frac{\Gamma\left(\frac{v+1}{2}\right)}{\Gamma\left(\frac{v}{2}\right)\Gamma\left(\frac{v}{2}\right)} \left(\frac{h}{v}\right)^{v/2} \left[1 + \frac{h}{v}(x-\theta)^2\right]^{-\frac{(v+1)}{2}} \quad (A.1)$$

$-\infty < x < \infty$   
 $-\infty < \theta < \infty$   
 $0 < h < \infty$   
 $v > 0$

### 2. LA DISTRIBUCION GAMMA ( $\alpha, \beta$ )

Una variable aleatoria (v.a.) se distribuye como una Gamma ( $\alpha, \beta$ ) si se expresa como:

$$p(x/\alpha, \beta) = \frac{x^{\beta-1} e^{-x/\alpha}}{\Gamma(\beta) \alpha^\beta} \quad (A.2)$$

$0 < x < \infty$   
 $\alpha, \beta > 0$   
 $\alpha$ : parámetro de escala

### 3. LA DISTRIBUCION $\chi^2$ CUADRADA.

Esta distribución es un caso particular de una Gamma ( $\alpha, \beta$ ), pues si se substituye en la expresión (A.2)  $\alpha=2$  y  $\beta=\frac{v}{2}$ , se obtiene la cuadrada.

$$p(x/v) = \frac{x^{v/2-1} e^{-x/2}}{2^{v/2} \Gamma\left(\frac{v}{2}\right)} \quad (A.3)$$

$0 < x < \infty$   
 $v > 0$ , indica el número de grados de libertad.

### 4. GAMMA INVERTIDA.

Esta distribución es obtenida de (A.2) haciendo un cambio de variable adecuado, dado por  $y^2=1/x$ , de donde

$$p(y/\alpha, \beta) = \frac{2}{\Gamma(\alpha) \alpha^\beta} y^{2\beta+1} e^{-\frac{1}{\alpha y^2}} \quad (A.4)$$

$0 < y < \infty$   
 $\alpha, \beta > 0$

Otra expresión, de la Gamma invertida es tomar  $\sigma = y$ ,  $\rho = \frac{y}{2}$  y  $\alpha = \frac{2}{\sigma^2}$ , con lo cual

$$P(\sigma/v, S) = \frac{2}{\Gamma(\frac{\alpha}{2}) (\frac{2}{\sigma^2})^{\frac{\alpha}{2}} y^{\alpha+1}} e^{-\frac{y^2}{2\sigma^2}} \quad \begin{matrix} 0 < \sigma < \infty \\ v, S > 0 \end{matrix} \quad (A.5)$$

5. ZELLNER DEFINE UNA FUNCIÓN DE PERDIDA SIMÉTRICA COMO AQUELLA QUE SATISFAGA LO SIGUIENTE

$$L(H_0, H_1) = L(H_1, H_0) \quad \text{y} \quad L(H_0, H_0) = L(H_1, H_1) = 0 \quad (A.5)$$

Donde  $H_0$  y  $H_1$  son hipótesis exclusivas y exhaustivas.

Para mayor información ver Zellner (1971, pp. 295-297)

6. LAS REGLAS DE JEFFREYS (ZELLNER, 1971, PP. 42-44) SIRVEN PARA SELECCIONAR UNA DISTRIBUCIÓN INICIAL QUE REPRESENTA POCO CONOCIMIENTO O IGNORANCIA. ESTO CONSISTE EN

- a) Si se cuenta con el parámetro  $\mu$  desconocido, la ignorancia se representa como

$$P(\mu) \propto \text{constante} \quad - \quad 0 < \mu < \infty$$

- b) Si  $\sigma$  es desconocida, una representación acerca de la ignorancia es:

$$P(\sigma) \propto \frac{1}{\sigma}$$

## BIBLIOGRAFIA

- ATKINSON, A.C. (1978), "Posterior probabilities for choosing a regression model". *Biometrika* 65, 39-48.
- BERNARDO, J.M. (1979), "Reference posterior distributions for Bayesian Inference", *Journal of the Royal Statistical Society*, B, 41, pp.113-147.
- BERNARDO, J.M. (1980), "Conceptos, Métodos y Fuentes de la Bioestadística", *Comunicación Interna* No.1.
- BROEMELING, L.D. (1985), "Bayesian analysis of linear models", Marcel Dekker, New York, pp.94-104.
- BROOKS, R.J. (1972), "A decision theoretic approach to optimal regression designs", *Biometrika*, 59, 563-71.
- BROOKS, R.J. (1974), "On the choice of an experiment for prediction in linear regression", *Biometrika*, 61, 303-311.
- BOX, G.E.P. & TIAO, G.C. (1973), "Bayesian inference in Statistical Analysis", Addison Wesley.
- CHATFIELD, C & COLLINS, A.J. (1980), *Introduction to Multivariate Analysis*", Science Paperbacks.
- CHATTERJEE, S & BERTRAM, P. (1977), *Regression analysis by example*", Wiley.
- DAVIS, W.W. (1979), "Approximate bayesian predictive distributions and model selection". *Journal of the American Statistical Association*, V.74, pp.312-317.



- DEGROOT, M. H. (1970), "Optimal Statistical Decisions", McGraw-Hill. New York.
- DRAPER, N. R. & SMITH, H. (1981), "Applied Regression Analysis". Second Edition, Wiley. New York.
- GEISSER, S. (1975), "The predictive sample reuse method with applications", *Journal of the American Statistical Association*, V. 70, pp. 320-328.
- GEISSER, S. & EDDY, W. F. (1979), "A predictive approach to model selection". *Journal of the American Statistical Association*, V. 74, pp. 153-168.
- GOLDSTEIN, M. (1976), "Bayesian Analysis of Regression Problems", *Biometrika*, 63, pp. 51-58.
- HAITOVSKY, Y. (1972), "Regression Estimation From Grouped Observations", No. 33, *Griffin's Statistical Monographs & Courses*. Alan Stuart. D. Sc. London.
- JEFFREYS, H. (1961), "Theory of Probability (3rd ed)". Oxford Clarendon.
- LINDLEY, D. V. (1968). "The Choice of Variables in Multiple Regression", *Journal of the Royal Statistical Society*, B. V. 30, pp. 31-66.
- LINDLEY, D. V. (1977), "Principios de la Teoría de la Decisión". Ed. Vicen-Vives. Barcelona.
- MONTGOMERY, D. C. & PECK, E. A. (1982), "Introduction to Linear Regression Analysis", John Wiley.
- PERICHI, L. R. (1984), "An alternative to the Standard Bayesian Procedure For Discrimination Between Normal Linear Models". *Biometrika*, V. 71, pp. 575-586.
- STONE, M. (1974), "Cross-Validatory Choice & Assessment of Statistical Predictions", *Journal of The Royal Statistical Society*, B, 36, pp. 111-147.

STONE, M. (1976), *An Asymptotic Equivalence of Choice of Model by cross-validation and Akaike's criterion*, *Journal of the Royal Statistical Society, Soc. B*, 39, pp. 44-47.

WEISBERG, S. (1980), *"Applied linear Regression"*. John Wiley & Sons, Inc.

Zellner, A (1971), *"An Introduction to Bayesian Inference in Econometrics"*, Wiley, New York.