Universidad Nacional de Ingeniería

PROGRAMA ACADEMICO ESCUELA DE GRADUADOS



"ANALISIS Y SELECCION DE VARIABLES EN REGRESION LINEAL "

TESIS

Para optar el Grado de Magister en Ciencias, Mención en:

MATEMATICAS APLICADAS

Presentada por

CARLOS A. BAZAN CABANILLAS

LIMA . PERU . 1980

INTRODUCCION

La técnica de la Regresión es, sin lugar a dudas, uno de los puntos de mayor aplicación en Estadística, sobre todo la regresión de tipo lineal, ya que, en problemas linea les es más viable llegar a resultados tangibles con menor es fuerzo. Sin embargo, el problema de la regresión lineal continúa siendo un problema abierto, siendo la selección óptima de las variables de regresión y la multicolinealidad de los detos dos espectos importantes en actual discusión.

Una razón fundamental para que este problema no ha ya eido resuelto, es que no ha sido bien definido. Parece que no hay un único problema, sino que ceda uso particular que se da e la regresión define un problema, y para cada ceso se hace necesaria una respuesta. Por otra parte, razones de orden económico, técnico y estadístico exigen la determinación del "major subconjunto de regresión". El estudio de estos espectos se ha visto impulsado últimamente por el apoyo que - brinde la ciencia de la computación.

El propósito del presente trabajo es proporcionar los conceptos fundamentales de la regresión lineal y los critarios para la selección de las variables en dicho modelo. Se estudiará el problema de la multicolinealidad de las variables de predicción y se expondrá la técnica de regresión por ariatas, con el objeto de superar las inconveniencias de la multicolinealidad.

El problema de la selección de variables consta de tres pasos: 1) El uso de algoritmos que proporciona información para el análisis, por ejemplo se podría usar regresión

por pasos.

- 2) El ueo de criterios pera analizar y seleccionar las variables, por ejemplo, determinar cuando se debe detener la regresión por pasos, y
- 3) La estimación de los coeficientes de la ecuación de regresión, por ejemplo por eliminación Gaussiana.

En la solución de un problema concreto, la obtención do "datos buenos" presupone la homogeneidad de las varianzas, la ausencia de datos erráticos, la preferencia por la ortogo nalidad de las columnas de la matriz de diseño de experimen - tos, etc. Un problema seriojes la ausencia de ortogonalidad (multicolinealidad) lo cual origina problemas de estimación de parámetros, de selección de variables y de interpretación de resultados empíricos. Por lo tanto, se hace necesario la de terminación de métodos robustos a la multicolinealidad, en eg te sentido la regresión por aristas trata de superar este problema.

Los dos únicos tipos de ajusta de curvas que se diginate de curvas que se diginate en serán; el de mínimos cuadrados, que se usa en los Capítulos Nos. 1, 2, 3 y 4; y el de aristas que se usa en el Capítulo No.5.

En el Capítulo No.1 se presente la notación y los conceptos fundamentales; en el Cap.2 se presentan los algoritmos computacionales; el Cap.3 nos proporciona los criterios para la salección de variablas; en el Cap.4 se presenta el problema de multicolinealidad; y, en el Cap.5 se ofrece la alternativa de regresión por aristas, con el objeto de superar los inconvenientes de la multicolinealidad.

Carlos Bazán C.

Lima, abril de 1978

CONTENIDO

		pp.
Capítulo No. 1	Notación y conceptos básicos.	1
Capítulo No. 2	Técnicas de selección de variables.	14
Capítulo No. 3	Critorios de selección	21
Capítulo No. 4	Multicolinealidad	35
Capítulo No. 5	Regresión por aristas	48
	Resumen	56
	Bibliografía	58

CAPITULO No 1

NOTACION Y CONCEPTOS BASICOS

1.1 Noteción e hipótesis.

Asumamos que existe una relación lineal entre t $v_{\underline{a}}$ riables: $x' = (x_1, x_2, ..., x_t)$, una variable aleatoria de perturbación e, y una variable dependiente y = y(x, e).

Si tenemos n (n)t + 1) observaciones de y, a la vez que de x, podemos escribir

$$y_i = \beta_0 + \sum_{j=1}^{t} x_{ij} \beta_j + e_j, 1 \le i \le n.$$
 [1.1]

Los coeficientes β y los parámetros de la distrib<u>u</u> ción de ℓ son desconocidos, y el problema consiste en obtener sus estimaciones. En forme matricial se tiene

$$Y = X\beta + \theta , \qquad (1.2)$$

donde :

Y es un n vector columna de observaciones da y_{τ}

X es la n \times (t + 1) matriz de diseño de experimon tos definida según (1.1) Los elementos de la primera columna de \times son iguales a 1.

 β es un (t + 1) vector columne, no conocido apriori; y

€ os un n vector columna aleatorio no observado, d<u>e</u> finido según (1.1).

Asumamos también las hipótesis (1.3)

- 1.3.a) E(e) = 0,
- 1.3.b) V(e) = E(ee') = 6³In,
- 1.3.c) X es une metriz de rengo t + 1
- 1.3.d) En $x_1, x_2, ..., x_t$ están incluidas todas las variables importantes en la determinación de y mediante (1.2), pero

también pueden estar incluidas otras no importantes, en la práctica esta suposición debe ser tomada con cautela ya que as frecuenta que no se disponga de datos para al guna variable importante, esta problema se discuta en - la sección 1.2.

1.3.e En el proceso de selección de variables se eliminan rtérminos, esto equivale a fijar rcoeficientes β con el valor caro.

El término asociado a β_0 no será eliminado. Sin pérdida de ganeralidad, podemos suponer que se eliminan las y variables $x_{t-Y_{+1}}, \ldots, x_{t}, \quad \text{Sean p = t + 1 - r, los términos que se retienen.}$

- Hagamos dos convenciones:
- 1- Los parámetros estadísticos y las cantidades relacionadas con los modelos restringidos de p y r términos, serán indica dos por los subíndicos p y r respectivamente; y los relacionados con ol "modelo completo" de t + 1 términos no llevarán subíndice. Luego, particionando convenientemente,(1.2) se escribe

$$Y = X_{p} \beta_{r} + X_{r} \beta_{r} + \theta$$
 (1.4)

2- Debido a que en nuestro medio, la bibliografía disponibla está, en su mayoría, en el idioma Inglés, para facilidad de los lectores, cuando se usan siglas éstas se mantendran en su forma original, por ejemplo, 833 corresponde a "Residual sum of squares" (suma de los cuadrados de los residuales) del mode lo de p términos.

<u>Definición 1.1.</u> Entre todos los modelos de p términos, se llamará el "mejor" modelo de temaño p a aquel que terga mínimo ASS_p. Se anfatiza que el adjetivo "mejor" indica estrictamente mínimo field, y quo realmente el modelo podría no ser el mejor para otro uso que se le destine. Además, la definición de mejor es aplicada solo a la muestra en uso y esto no implica — que la relación se mantiene necesariamente para la población.

1.2 Consecuencias de especificar el modelo con un número de variables diferente el adecuado

Mazones de orden económico, técnico e inclusiva estadístico (por ejemplo, variabilidad de los parámetros estimados y de la predicción) hacen recomendable la eliminación de al gunas variables. En esta sección, bajo la suposición de linealidad del modelo, se hace una revisión de las consecuencias de una especificación incorrecta del modelo, ya sea por la eliminación de variables importante y/o por la retención de variables superfluas.

Estimación de los parámetros 6 y 62

Asumamos que el modelo "completo" de t + 1 térmi - nos es correcto, y sea

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_{P} \\ \hat{\beta}_{T} \end{bmatrix}$$

El estimador mínimo - cuadrático de β pera el modelo de t+1 términos. Entonces, el rango de X=t+1, se tiene que

$$\widehat{\beta} = (x'x)^{-1}x'Y, \qquad (1.5.a)$$

$$E(\hat{\beta}) = \beta \tag{1.5.6}$$

y

$$v(\hat{\mathbf{g}}) = e^{2}(x'x)^{-1}$$
. (1.5.c)

3e demucetra fácilmente que β co el cetimador lineal inecegado de mínima varianza de β. 3ec

$$\hat{Y} = X\hat{\beta}, \qquad (1.5.a)$$

la predicción de Y en X, definemos el vector de residuos ξ = Y - \hat{Y} , (1.5.b)

luego un estimador muy apropiado para Ǽ es

$$\hat{\delta}^{2} = \frac{\varepsilon' \mathcal{E}}{\text{gredos do liberted de } \left(z \text{ df}\right)}$$

$$= \frac{(Y - \hat{Y})'(Y - \hat{Y})}{\text{df}} = \frac{Y'Y - \hat{Y}'Y}{\text{df}}$$

$$= \frac{Y'(I - X(X'X)^{-1}X')Y}{\text{df}} = \frac{e'(I - X(X'X)^{-1}X')e}{\text{df}}$$

Como $M \leq I - X(X'X)^{-1}X'$ se idempotente, se tiene que su rango es iguel e su traza = n - t - 1 = dF y según (1.3.a) y (1.3.b)

$$E(6^2) = \frac{\text{traze M}}{D - t - 1} = 6^2$$
 (1.5.d)

por lo tento, 6^2 es insesgado y si $C \rightarrow N(0, 6^2)$, el corolario 4.7.1, Graybill(1961) nos indica que $(n - t - 1)\hat{6}^2 \rightarrow \chi^2_{n-t-1}$.

En el modelo restringido, si llamamos $\widetilde{oldsymbol{eta}}_{
ho}$ el estimador mínimo-cuadrático de $oldsymbol{eta}_{
ho}$, sé tiene

$$\vec{\beta}_{p} = (x_{p}^{2}x_{p}^{2})^{-1}x_{p}^{2}Y_{p}$$
 (1.7.e)

$$E(\hat{\beta}_{0}) = \beta_{0} + A \beta_{F}, \qquad (1.7.6)$$

dande

$$A = (x_{p}^{*}x_{p}^{*})^{-1}x_{p}^{*}x_{p}^{*}$$
(1. 8)

$$V(\hat{\beta}_{p}) = \delta^{2}(x_{p}^{2}x_{p}^{2})^{-1}$$
 (1.7.0)

Análogamento el caso anterior, un ostimador de 6º es

$$\tilde{e}^{2} = \frac{(Y-\tilde{Y})'(Y-\tilde{Y})}{dF} = \frac{y'(I_{e}X_{p}(X_{p}^{*}X_{p}^{*})^{-1}X_{p}^{*})Y}{dF}$$
 (1.9.a)

* Si Y tiene distribución $N(\mu, 6^2I)$, entonces $\frac{Y'AY}{6^2}$ tiene distribución $\chi^2(K,\lambda)$, donde $\chi = \frac{\mu'A\mu}{26^2}$, si y solo A es una matriz idempotente de rango K. Además $\varepsilon(\chi^2(K,\lambda)) = K + 2\lambda$.

y, si $e \longrightarrow N(0, e^2I)$, por sl corolario 4.7.1 de Graybill(1961) se tieno que

$$\frac{\partial^{2}(n-p)}{\partial^{2}} \longrightarrow \chi^{2}(n-p,\lambda) \text{ con } \lambda = \frac{\beta_{r}\chi_{r}^{\prime}(I-\chi_{r}(\chi_{r}^{\prime}\chi_{r}))^{-1}\chi_{r}^{\prime})\chi_{r}^{\beta_{r}}}{2\sigma^{2}}$$

Por lo tento

$$E(\hat{\mathbf{6}}^{2}) = G^{2} + \frac{\beta_{r}^{'} X_{r}^{*} (\mathbf{I} - X_{p} (X_{p}^{*} X_{p})^{-1} X_{p}^{*}) X_{r} \beta_{r}}{D - P}$$
(1.8.6)

El cetedístico que mide la variabilidad do $\hat{\beta}_p$ con respecto a su valor esperado $E(\hat{\beta}_p)$ es $V(\hat{\beta}_p)$ y el estadístico que mide la variabilidad de $\hat{\beta}_p$ con respecto a su valor real esperado β_p es la media de los cuadrados de los erroras (mean equare error (MSE($\hat{\beta}_p$)).

HSE(
$$\hat{\beta}_{p}$$
) = E($\hat{\beta}_{p}$ - β_{p})($\hat{\beta}_{p}$ - β_{p})' = E($\hat{\beta}_{p}$ -E($\hat{\beta}_{p}$)+A β_{r})($\hat{\beta}_{p}$ -E($\hat{\beta}_{p}$)+A β_{r})'

=E(($\hat{\beta}_{p}$ -E($\hat{\beta}_{p}$)($\hat{\beta}_{p}$ -E($\hat{\beta}_{p}$))') + 2E(($\hat{\beta}_{p}$ -E($\hat{\beta}_{p}$) β_{r} A') + A β_{r} B'A'

= V($\hat{\beta}_{p}$) + A β_{r} B'A'. (1.10)

Fradicción y estimación

Sea x un vector fijo de ingreso de t+1 componentes, sea la variable dependiente:

$$Y(x) \equiv y$$
. (1.11.a)

Según al modelo lineel sa tiano

$$E(y) = \tilde{y} = x^{i}\beta, \qquad (1.11.b)$$

la estimación $\widehat{\hat{\mathbf{y}}}$ de $\widehat{\mathbf{y}}$, la esperanza y la varianza de predicción de $\widehat{\hat{\mathbf{y}}}$ están dadas por

$$\hat{\vec{y}} = x'\hat{\beta} , \qquad (1.11.c)$$

$$\vec{z}(\hat{\vec{y}}) = \mathbf{x'}\boldsymbol{\beta} \tag{1.11.d}$$

У

$$V(\hat{S}) = x'(X'X)^{-1} \times 6^2$$
 (1.11.e)

Para la predicción simple, so tieno

$$\widehat{Y}(x) \equiv \widehat{y}$$
 (1.12.e)

como y = ỹ+ệ,

$$\hat{y} = \hat{y} + \hat{e} = \hat{y} + \hat{e} , \qquad (1.12.6)$$

con ŷ y ê independientes, luego

$$E(\hat{y}) = x'\beta \qquad (1.12.0)$$

У

$$VP(\hat{y}) = \kappa^2 (1 + x'(X'X)^{-1} x). \qquad (1.12.d)$$

Analogamente, an al modelo restringido, para un \mathbf{x}_{p} de p variables

$$Y_{p}(x_{p}) = y_{p}$$
 (1.13.2)

У

$$E(y_p) \equiv \overline{y}_p = x_p^* \beta_p . \qquad (1.13.6)$$

Para una estimación $\overset{\smile}{y}$ ds \overline{y}_p , se tiens

$$\tilde{\vec{y}}_{p} = x_{p}^{\prime} \hat{\beta}_{p}^{\prime} , \qquad (1.13.c)$$

$$E(\tilde{\vec{y}}) = x_0^2 \beta_D + x_0^2 A \beta_D \qquad (1.13.d)$$

У

$$V(\tilde{y}) = s^2 x_p^2 (x_p^2 x_p)^{-1} x_p$$
 (1.13.e)

El estadístico que mide la variabilidad de $\tilde{\vec{y}}$ con respecto a su media as $\forall (\tilde{\vec{y}})$, y el estadístico que mide la variabilidad de $\tilde{\vec{y}}$ con respecto al valor real esperado $\tilde{y} = x^{2}\beta$, es la media del cuadrado del error de estimación.

$$MSE(\widetilde{y}) = E(\widetilde{y} - \widetilde{y})^{2}. \qquad (1.13.f)$$

En el Teorema 1.1 se verá que

$$MSE(\tilde{y}) = V(\tilde{y}) - (x_{\tilde{p}}A\beta_{r}-x_{\tilde{p}}\beta_{r})^{2}. \qquad (1.13.g)$$

Y, para predicción simple de $Y_p(x_p) = \widetilde{y}$, se tiene

$$\tilde{y} = \tilde{y}_p + \theta$$
, (1.14.a)

$$E(\tilde{y}) = x_p^2 \beta_p + x_p^2 A \beta_r , \qquad (1.14,b)$$

$$VP(\tilde{y}) = e^{2(1+x_{P}^{2}(X_{P}^{2}X_{P})^{-1}x_{P}^{2})$$
 (1.14.c)

У

$$MSEF(\widetilde{y}) = E(\widetilde{y} - y)^{2} . \qquad (1.14.d)$$

En el Teorema 1.1 se verá que

$$MSEF(\tilde{y}) = VF(\tilde{y}) - (x_0^2 A \beta_D - x_0^2 \beta_D)^2 . \qquad (1.14.6)$$

Con los conceptos anteriores probaremos seis propiadedes muy importantes:

Toorema 1.1.

- 1) $\tilde{\beta}$ as sesgado, excepto cuando $X_{\tilde{\rho}}^{1}X_{r}\beta_{r}=0$ Y $\tilde{\zeta}^{2}$ as sesgado, excepto cuando $\beta_{r}=0$.
- 2) $V(\widehat{\beta}_p)$ $V(\widehat{\beta}_p)$ es semidefinida positiva. Esto indica que los estimadores de las componentes de $\widehat{\beta}_p$ dados por $\widehat{\beta}_p$ son generalmente más variables que los dedos por $\widehat{\beta}_p$, ya que $V(\widehat{\beta}_i) > V(\widehat{\beta}_i)$.
- 3) 9i la matriz $V(\beta_r) = \beta_r \beta_r'$ sa semidefinida positiva, se tiena que

 $\begin{array}{l} \text{V}(\widehat{\beta_p}) \text{ - MSE}(\widehat{\beta_p}) \text{ as semidefinide positive, for lo tento,} \\ \text{V}(\widehat{\beta_i}) \text{ - MSE}(\widehat{\beta_{pi}}) \geqslant 0 \text{ .} \end{array}$

4)
$$\text{MSE}(\tilde{\vec{y}}) = V(\tilde{\vec{y}}) - (x_{\vec{p}}^* A \beta_{p} - x_{\vec{p}}^* \beta_{p}^*)^2$$
 [1.13.g]
$$y$$

$$\text{MSEP}(\tilde{\vec{y}}) = VP(\tilde{\vec{y}}) - (x_{\vec{p}}^* A \beta_{p} - x_{\vec{p}}^* \beta_{p}^*)^2$$
 [1.14.e]

- 5) VP(ŷ) > VP(ÿ) .
- 6) Si, $V(\widehat{\beta}_r) = \beta_r \beta_r'$ es _ _ _ _ _ _ semidefinide positive, entonces $VP(\widehat{y}) \geqslant MSEP(\widehat{y})$.

<u>Prueba.</u> En toda esta prueba asumiremos la notación siguiento:

$$X = (X_p, X_r) ,$$

entonces,

$$X'X = \begin{bmatrix} X_{D}^{1}X_{D} & X_{D}^{2}X_{D} \\ X_{D}^{1}X_{D} & X_{D}^{2}X_{D} \end{bmatrix} \in C \equiv \begin{bmatrix} C_{DD} & C_{DD} \\ C_{DD} & C_{DD} \end{bmatrix} . \quad (1.15.e)$$

Sea

$$(x'x)^{-1} \equiv B \equiv \begin{bmatrix} B_{pp} & B_{pr} \\ B_{rp} & B_{rr} \end{bmatrix} . \tag{1.15.b}$$

Las matrices 8 y C son definides positivas. Con el objeto de su utilización posterior, hagamos el siguiente cálculo:

$$B_{pr} B_{rr}^{-1} = B_{pr}(C_{rr} - C_{rp}C_{pp}^{-1}C_{pr})$$

$$= B_{pr}C_{rr} - B_{pr}C_{rp}C_{pp}^{-1}C_{pr}$$

$$= - B_{pp}C_{pr} - (I - B_{pp}C_{pp})C_{pp}^{-1}C_{pr}$$

$$= - C_{pp}^{-1}C_{pr}$$

$$= - (X_{p}^{*}X_{p}^{*})^{-1}X_{p}^{*}X_{r} .$$

Según (1.8) so tiene

$$B_{\rm DF}B_{\rm DF}^{-1} = -A$$
 (1.15.c)

Perto 1. Según (1.7.6), $E(\vec{\beta}_p) = \hat{\beta}_p + (x_p'x_p)^{-1}x_p'x_p \hat{\beta}_p$, luago $\vec{\beta}_p$ as sasgado, excepto cuando $x_p'x_p\beta_p = B$.

De (1.15.a) y (1.15.c) se obtiene

$$B_{rr}^{-1} = X_{r}^{*}(I - X_{r}(X_{r}^{*}X_{r})^{-1}X_{r}^{*})X_{r}$$

la cual también es definida positiva, aplicando (1.9.b)

$$\epsilon(\tilde{c}^2) = c^2 + \frac{\beta r \mathbf{6}^{-1} \beta r}{\beta - 2}$$
,

Finalmento, se tiono que $\widetilde{\delta}^2$ os sesgada, excepto cuando $eta_r=0$.

Parts 2. Según (1.15.a) y (1.15.b) vemos que

$$a_{pp}^{-1} = (x_p^*)_p^{-1} = a_{pp}^{-1} - a_{pp}^{-1} a_{pp}^{-1}$$
.

Como $V(\hat{\beta}) = (x'x)^{-1} \delta^2$, se tisne que $V(\hat{\beta}_p) = B_{pp} \delta^2$ y

 $v(\tilde{\beta}_p) = (x_p^2 x_p)^{-1} \sigma^2$, luago:

$$V(\hat{\beta}_{p}) - V(\hat{\beta}_{p}) = B_{pr}B_{rr}^{-1}B_{rp}\sigma^{2}$$
. (1.15.d)

Veamos que esta matriz es semidefinida positiva. Sea $Z_{
m p}$ un vector cualquiera de p componentes, entonces

$$Z_p^* \theta_{pp} \theta_{pp}^{-1} \theta_{pp} \sigma^2 Z_p = \sigma^2 (\theta_{pp} Z_p) '\theta_{pp}^{-1} (\theta_{pp} Z_p) \geqslant 0 .$$

puesto que θ_{rr}^{-1} es definida positiva.

Si sa toma
$$Z_p^* = \{0, ..., 0, 1, 0, ..., 0\}$$
, se tiene

 $Z_{i}(V(\hat{\beta}_{p}) - V(\hat{\beta}_{p}))Z_{p} = V(\hat{\beta}_{i}) - V(\hat{\beta}_{i}) > 0$ luego, $V(\hat{\beta}_i) \ge V(\hat{\beta}_i)$.

Parts 3, Según (1.10) y (1.15,d) se tiene MSE $(\widetilde{\beta}_{D}) = V(\widetilde{\beta}_{D}) + A\beta_{D}\beta_{D}^{2}A$ $V(\widetilde{\beta}_{D}) = V(\widehat{\beta}_{D}) - S^{2}B_{D}B_{D}^{2}B_{D}^{2}$ por lo tanto $V(\hat{\beta}_D)$ - MSE $(\hat{\beta}_D)$ = $\delta^2 B_{DP} B_{PD}^{-1} B_{PD}$ - $A \beta_P \beta_P^* A^*$.

Sea Zo un vector qualquiera de p componentes, entonces

$$Z_{p}^{*}(V(\hat{\beta}_{p}) - MDE(\hat{\beta}_{p}))Z_{p} = Z_{p}^{*}(\sigma^{2}\theta_{p}rB_{r}^{-1}\theta_{rp}-A\beta_{r}\beta_{r}^{*}A^{*})Z_{p}$$

Según (1.15.c) se tiene

$$z_{i}(v(\hat{\beta}_{p})-\text{MSE}(\hat{\beta}_{p}))z_{p} = z_{i}(s^{2}\theta_{p}-\theta_{r}^{-1}\theta_{rp}-\theta_{p}-\theta_{r}^{-1}\theta_{r}\theta_{r}^{-1}\theta_{rp})z_{p}$$

$$= z_{i}^{2}\theta_{p}-\theta_{r}^{-1}(\theta_{r}-\sigma^{2}-\theta_{r}\theta_{r}^{2})\theta_{r}^{-1}\theta_{rp}z_{p}$$

$$= (\theta_{r}^{-1}\theta_{rp}z_{p})'(v(\hat{\beta}_{r})-\beta_{r}\theta_{r}^{2})(\theta_{r}^{-1}\theta_{rp}z_{p})$$

$$\geq 0.$$

ya que por hipotesia, $V(\widehat{\beta}_r)$ - $\beta_r \beta_r^*$ as somidefinide positiva. Si se tome $Z_{p}^{2} = \{0, ..., 0, 1, 0, ..., 0\}$, se tione

 $Z_{G}(V(\widehat{\beta}_{D}) - MSE(\widehat{\beta}_{D})) Z_{D} = V(\widehat{\beta}_{i}) - MSE(\widehat{\beta}_{i}) \geqslant 0, lusgo$ $\forall (\hat{\beta_i}) \geqslant MSE(\hat{\beta_i})$.

Perte 4. Probaremos solamente (1.14.e), la pruebe de (1.13.g) es similar.

Calculemos directamenta NAEP.

мэєр(ў) =
$$E(y-\vec{y})^2 = E((y-\vec{y})(y-\vec{y}))$$

= $E((x_p^2\beta_p + x_p^2\beta_p + e - x_p^2\beta_p)(x_p^2\beta_p + x_p^2\beta_p + e - x_p^2\beta_p))$
= $E(x_p^2(\beta_p - \hat{\beta}_p)(\beta_p - \hat{\beta}_p))(x_p + e + e + e)$
+ $(x_p^2\beta_p + e)(x_p^2\beta_p + e)$
= $x_p^2E((\beta_p - \hat{\beta}_p)(\beta_p - \hat{\beta}_p))(x_p + e + e)$
+ $e^2((\beta_p - \hat{\beta}_p)(\beta_p - \hat{\beta}_p))(x_p + e)$
+ $e^2((\beta_p - \hat{\beta}_p)(\beta_p - \hat{\beta}_p))(x_p + e)$

$$= x_{p}^{2}MGE(\widetilde{\beta}_{p})x_{p} - 2x_{p}^{2}A\beta_{p}\beta_{r}^{2}x_{p} + 0 + x_{p}^{2}\beta_{p}\beta_{r}^{2}x_{p} + 0 + 6^{2}$$

$$= x_{p}^{2}(V(\widetilde{\beta}_{p}^{2}) + A\beta_{p}\beta_{r}^{2}A^{2})x_{p} + 6^{2} - 2x_{p}^{2}A\beta_{p}\beta_{r}^{2}x_{p} + x_{p}^{2}\beta_{p}\beta_{r}^{2}x_{p}$$

$$= VP(\widetilde{y}) + (x_{p}^{2}A\beta_{p} - x_{p}^{2}\beta_{p})^{2} .$$

Parte 5, Probemos antes un Lema.

Lame. Sea M une matriz k x k definide positive, entonces

$$\bar{M} = \begin{bmatrix} M^{-1} & I \\ I & M \end{bmatrix}$$

os une metriz 2k x 2k (definide positive.

<u>Prueba.</u> See P una matriz ortogonal que diagonaliza a M, luego

$$P'MP = D = \begin{bmatrix} d_{1}_{d_{2}} & 0 \\ 0 & \cdots \\ d_{k} \end{bmatrix},$$

los d. > O ya qua M es definida positiva. sea

$$\vec{F} = \begin{bmatrix} \vec{P} & 0 \\ 0 & \vec{P} \end{bmatrix} ,$$

P es ortogonal ya que P lo es, luego

$$\vec{P}^{\dagger}\vec{M}\vec{P} = \begin{pmatrix} 0^{-1} & I \\ I & D \end{pmatrix} = \begin{pmatrix} 1/d_1 & 0 & 1 & 0 \\ 0 & 1/d_k & 0 & 1 \\ 1 & 0 & d_1 & 0 \\ 0 & 0 & d_k \end{pmatrix} ,$$

esta metriz tiena los determinantes de los k primeros menores principales mayores que cero, y los k siguientes iguales a cero, por lo tento, $\vec{P}'\vec{M}\vec{P}$ es semidefinida positiva. Veamos que \vec{M} es semidefinida positiva; sea Z un vector cualquiera de Zk componentes, sea $V = \vec{P}'Z$. Luego, $Z'\vec{M}Z = V'(\vec{P}'\vec{M}\vec{P})V \geqslant 0$, por lo tento, \vec{M} es semidefinida positiva.

Prusba do la ponte 5.

set $x'=(x_p^*,x_n^*)$ un vector de ingreso, entonces, $\hat{y}=x_p^*\hat{\beta} \quad y \quad \hat{y}=x_p^*\hat{\beta} \quad , \ \, \text{de donde}$

$$VP(\hat{y}) = \sigma^{2}(x'(X'X)^{-1}x - x'_{p}(X'_{p}X_{p})^{-1}x_{p})$$

$$= \delta^{2}(x'_{p}B_{p}B_{p}^{-1}B_{p}x_{p} + x'_{p}B_{p}x_{p} + 2x'_{p}B_{p}x_{p}$$
(1.15.d)

Definamos Z' \in $\{x_p^{\prime}B_{pr}, x_r^{\prime}\}$, Z tieno 2r componentes; reemplezando en la expresión enterior, se tiene

$$VP(\hat{y}) - VP(\hat{y}) = Z' \begin{pmatrix} B_{rr}^{-1} & I \\ I & B_{rr} \end{pmatrix} Z$$
,

Aplicando el Lama anterior se tiene $VP(\hat{y}) \geqslant VP(\hat{y})$.

Parto B. Frobemos entes un Leme.

<u>Lema 1.2.</u> Sea M una matriz k x k semidefinide positiva entonces

$$\widetilde{M} = \begin{bmatrix} M & M \\ M & M \end{bmatrix}$$

es una matriz 2k x 2k somidofinida positiva.

<u>Prueba.</u> See P una matriz ontogonal que diagonaliza e M. luego

$$F'MF = D = \begin{bmatrix} d_1 & 0 \\ & d_2 & 0 \\ 0 & & d_k \end{bmatrix} ,$$

con $d_i \geqslant 0$ ya qua M es semidafinida positiva. Sea

$$\tilde{P} = \begin{bmatrix} P & 0 \\ 0 & P \end{bmatrix} ,$$

P es ortogonal ya que P lo es, luego

$$\vec{F}, \vec{M}\vec{F} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} d_1 & 0 & d_1 & 0 \\ 0 & d_k & 0 & d_k \\ d_1 & d_k & d_1 & d_k \\ \vdots & \vdots & \vdots & \vdots \\ 0 & d_k & 0 & d_k \end{pmatrix},$$

esta matriz tiene los determinantes de las K primeros menores

principales mayores o iguales a coro, y los k siguientes iguales a cero, por lo tento P'MP es semidefinida positiva.

Veamos que \widetilde{M} es semidefinida positiva. Sea Z un vector cualquiera de 2k componentes, Sea $V = \widetilde{F}'Z$, luego $Z'\widetilde{M}Z = V'(\widetilde{F}'\widetilde{MF})V \geqslant 0$, por lo tento \widetilde{M} es semidefinida positiva.

Fruebs de la parte 5. Sea x' = $\{x_p^i, x_p^i\}$ un vector de ingreso, entonces $\hat{y} = x'\hat{\beta}$ y $\hat{y} = x'\hat{\beta}_p$. Según (1.15.d), (1.14.c) y (1.15.c) se tiene

$$VF(\hat{y}) = e^{2}(x_{\hat{p}}B_{pr}B_{rp}^{-1}B_{rp}x_{p} + x_{\hat{p}}B_{rr}x_{r} + 2x_{\hat{p}}B_{pr}x_{r}),$$

$$VF(\hat{y}) = VF(\hat{y}) + (x_{\hat{p}}^{\prime}A\beta_{r} - x_{r}^{\prime}\beta_{r})^{2}$$

У

con lo cual se consigue

$$VP(\hat{y}) - MSEP(\hat{y}) = \sigma^{2}(x_{p}^{\prime}B_{p}B_{r}^{-1}B_{rp}x_{p} + x_{p}^{\prime}B_{rr}x_{r} + 2x_{p}^{\prime}B_{pr}x_{r})$$

$$- (x_{p}^{\prime}A\beta_{r} - x_{p}^{\prime}\beta_{r})^{2} \qquad (1.15.c)$$

$$= \sigma^{2}(x_{p}^{\prime}B_{pr}B_{r}^{-1}B_{rp}x_{p} + x_{p}^{\prime}B_{rr}x_{r} + 2x_{p}^{\prime}B_{pr}x_{r})$$

$$- x_{p}^{\prime}B_{pr}B_{r}^{-1}\beta_{r}\beta_{r}^{\prime}B_{rp}^{\prime}x_{p} - x_{p}^{\prime}\beta_{r}\beta_{r}^{\prime}x_{r}$$

$$- 2x_{p}^{\prime}B_{pr}B_{r}^{-1}\beta_{r}\beta_{r}^{\prime}x_{r} .$$

Offinamos Z' $= (x_p^* \theta_{pr} \theta_{rr}^{-1}, x_r^*)$, Z tiens 2r componentes; resmplazando en la expresión enterior so tiene

$$VP(\hat{y}) - MBEP(\hat{y}) = Z' \begin{pmatrix} c^{2p}rr - \beta r \beta \hat{r} & c^{2p}rr - \beta r \beta \hat{r} \\ c^{2p}rr - \beta r \beta \hat{r} & c^{2p}rr - \beta r \beta \hat{r} \end{pmatrix} Z,$$

como $V(\hat{\beta}_r) = e^{\hat{Z}_B}_{rr}$ y por hipotesis $V(\hat{\beta}_r) - \hat{\beta}_r \hat{\beta}_r'$ es semidefida positiva; aplicando el lema anterior se tiene $VF(\hat{y}) \geqslant MSEP(\hat{y})$, que es lo que queríamos prober.

OBSERVACIONES

- Las propiedades (2) y (5) recomiendan la eliminación de variables.
- Las propiedades (3) y (6) dan una condición bajo la cual la precisión del modelo restringido es mejor que la del modelo completo no obstante el sesgo.
- Si las variables x_r son extrañas, esto es β_r = 0, las propiedades (2) y (5) indican una pérdida de precisión en la estimación y predicción de las variables incluidas.

CAPITULO No 2#

TECNICAS DE SELECCION DE VARIABLES

En esta trabajo con sideramos el problema general que se presenta al tratar de determinar les rolaciones entre les variables de predicción x,y sus roles, tomadas esparadas o en conjunto al describir la respuesta y. Nuestro objetivo es la selección de un subconjunto de variables de predicción para la ecuación final.

rara proporcionar información, se hace necesario cons<u>i</u>

derar ajustes con diferentes combinaciones de variables de

predicción. Si el número de variables de predicción t, as peque

ño, se pueden considerar las 2^t combinaciones; pero si t es gran

de, tal posibilidad resulta enticconómica.

Describiremos los resgos principales de algunos métodos computacionales, proferentementa de mínimos cuadrados, aunque tembién mencionaremos el método de regresión por aristas
(ridge regresion). En el capítulo No3 se presenta una discusión de la forma adecuada de interpretar los resultados.

2.1.-Todas las regresiones posibles

Si el número t de variables de predicción no es muy grande, se puede considerar todas las 2º regresiones. En le actualidad hay un buen número de algoritmos que tratan este problema (ver Hocking 1976 pg.7). For ejemplo Morgan y Tatar trabajan con más de diez variables. La idea básica es la de realizar el cómputo en tal forma que dos conjuntos consecutivos difieran en una sola variable. Otra idea, un tanto distinta, desarrolleda por Newton y Spurrel(1967) es la que

*Una idea básica de les técnicas de pelección de variables se da en el Cap.No 6 de Draper(1956).

considera que la suma de los cuadrados debido a la regresión : $(\sum (\hat{Y}_i - Y_i)^2)$ es la suma de"elementos " básicos, y desarrollan un esquema para calcular estos elementos básicos, sin necesidad de evaluar todos los conjuntos.

La comparación de la eficiencia de estos algoritmos no es sencilla y las referencias no proporcionan información suficiente. Sin embargo, para el caso en que .la selección se ha ce en base a la suma de los cuadrados de los residuales, se informa que el segundo algoritmo de Furnival parece eficiente.

2.2.-Métodos de regresión por pagos

Debido el gran volumen de cálculos en la evaluación de todas las 2^t regresiones, se han propuesto otrus métodos en los que se evalúa un número menor de alternativas, ya sea por incremento o por eliminación de una variable en cada paso. A estos métodos se los llama de regresión por pasos y son la variación de dos ideas básicas llamadas selección hacia adelante (Forward selection: FS) y eliminación hacia atrás (Backward elimination: BE) Draper (1966) proporciona una exposición detallada. Vesmos los lineamientos generales de estos métodos:

Selección hacia adolunta: Este algoritmo empieza sin ninguna variable, y va adicioando una a una las variables, has_te que todas sa incluyan o hasta que se satisface un critario de finalización. Entre las variables hábiles para admisión, se eliga la que tiene el mayor coeficiento de correlación parcial con Y dado que han sido admitidas las otras variables. Esto es, la variable i as adicionada a la ecuación de p términos si:

$$F_{i} = \max_{j} F_{i} = \max_{j} \left(\frac{RSS_{p} - RSS_{p+jj}}{F_{hi}} \right) > F_{in}$$
 (2.1)

donde el subíndice p + j se refiere a las cantidades calculadas cuendo la variable j es adicionada al grupo de p variables.f. sigue dist.f(1,n -p - 1).La especificación del valor
fin determina una regle de finalización del algoritmo.
En la sección 3.4 se de un breve resumen de las reglas de finalización y los resultados de un estudio do simu_leción.

Eliminación hacia atrás : En este, caso so empieza conla cousción que incluye a todas las variables y luego se las va
aliminando una a una, La variable que se elimina es la que tiene el menor coeficiente de correlación percial, dedo que se man
tienan las otras variables. Esto es, la variable i es eliminada
de la ecuación de p términos si:

$$F_i = \min_j F_j = \min\left(\frac{RSSp - J - RSSp}{S^2}\right) < F_{aut} \qquad (2.2)$$

donde el subíndice p-j se refiere e la RSS calculade, cuendo la variable j se eliminada del grupo de p variables. f sigue dist. f(1, n - p). La especificación del valor f_{out} determina una regla de finalización del algoritmo.

A partir de estes dos ideas se deserrol)s otros com binaciones, siendo la más aceptada, la Jescrita por Efroymson -[ES]la cual es esoncialmente el proceso 75 en el que se incluye la posibilidad de oliminar una variable en cada pasa. Comunmenta, el proceso ES se la llama regresión por pasos.

Los procesos do regresión por pasos son criticados por no cumplir en general diversas condiciones recomendables, entre ellas:

- 1.-No se garantiza en forma absolute la obtención del"mejor" subconjunto de p variables.
- 2.-A menudo, el orden de inclusión y aliminación de variables no es el más adecuado .

3.-Las subrutinas típicas de cómputo proporcionan solamente la regresión de un solo conjunto de p variables, el cuel no es nacesariamente el mejor Además, como observara Gorman yToman(1958) el mejor subconjunto de p variables — no siempre es único.Por lo tento, es desemble obtener información de las mejores combinaciones de p variables para tomer una decisión final.Esta última condición es importente.

La idea intuitiva de estos procesos es que, para proble mas moderadamente bien definidos, el subconjunto de p variables elegido puede coincidir con el que se elige para todas las-regresiones, y el volumen de cálculos es relativamente pequeño. Un proceso ideal sería equel que para cada subconjunto de p variables de la mejor combinación y, además, otras combinaciones cercanas a la ideal, con un volumen de cálculos no muy elevado. En la siguiente sección se dan harramientas sobre este proceso.

2.3.<u>Subconjuntos óptimos</u>

Verios autores han desarrollado métodos que permiten determinar la bondad de las combinaciones de p variables.Le idea básica es la desarrollada por Hockin y Leslie(1967), que consista en lo siguiento: Supongamos que nuestro objetivo es obtener el mejor subconjunto de regresión que elimina a K vo. riables; para ello, ordenemos las variables de menor a mayor - según el coeficiente de regresión parcial al eliminar una variable dado que se mantienen todas las otras. Denotamos por - Q (i, j,...) a la suma de los cuadrados de los residuales cuamo do las variables $\frac{X_i}{i}, \frac{X_i}{i}, \dots$ son eliminadas, luego $\mathbf{Q}(i) \leq \mathbf{Q}(i,\dots)$. Supongamos que se consigue $\mathbf{Q}(1,2,\dots,k) \leq \mathbf{Q}(k+1)$, por la desigualdad anterior se tiene que $\mathbf{Q}(1,2,\dots,k) \leq \mathbf{Q}(k+1,\dots)$.

For 10 tento, les k variables que nos conviens eliminer son : $X_1, X_2, ..., X_k$. Si Q(1,2,...,k) > Q(k + 1) se evalúan otros — conjuntos.La extensión de esta idea es desarrollada e incorporada en un programa llemado SELECT, cuya eficiencia se mida por el número de evaluaciones que se requiero para obtenar un subconjunto óptimo que elimina a k variables, hay que conside rar que esta número de evaluaciones depende también de los datos particulares que se usan; una idea de la eficacia nos la da, el número de evaluaciones que se requiero para determinar la mejor combinación de variables para cada tamaño; para dos conjuntos particulares de datos de 15 y 26 variables, SELECT requirió de la evaluación de 1,465 y 3,546 subconjuntos de un to tel de $2^{15} = 32,768$ y $2^{25} = 678,108,364$ subconjuntos posibles respectivemente.

Otras alternativas en las que se considera la ob tención de los mejores subconjuntos para cada tamaño, y otras modalidades de evaluación se mencionan en Hocking(1976).

2.4.-Métodos subóptimos

Debido a la inseguridad de los métodos de ragresión por pesos y el alto volumen de cálculos que requieren los métodos de selección óptima, se desarrolla alternativas intermedias, entre ellas Barr y Goodnight(1971) en el Statistical — Analysis Sistem (SAS) Regression Program, proponen un esquema basado en el máximo cracimiento del estadístico R². Este proceso de se esencialmente la aplicación del proceso de Efroymson (ES) que so desarrolla de la forma siguiente: Si se tiene la mojor combinación de m-1 variables se incrementa — la variable que da el máximo cracimiento de R², y se realiza el intercambio de una variable incluida con una excluida,

No he encontrado un trabajojen el que se cotejen es tos métodos; pero en Hocking (1975) se mencione que verios usus - rios han reportado diferencias.

2.5.-Regresión por aristas (Ridge Regression)

Para problemas que involucran estadísticos de pre dicción ortogonales, Hoerl y Kennard (1970) sugieren el estimador sesgado de "Ridge"

$$B = B(K) = (x'x + K I)^{-1} x'Y,$$
 (2.3)

donde X ha sido estandarizado.La constante K se determina por inspección de la"traza" de arista,o sea el gráfico de B(K) ver sua K.En el Capítulo 5, desarrollaremos con mayor detalle las i deas del método de regresión per atistas, que es uno de los temas más importantes de este trabajo.

En el contexto do esta sección, nótese que aunque la regresión por aristas no está diseñada explicítamente para la eliminación do variables, hay inherente una eliminación de las variables cuyos coeficientes en(2.3) tienden rápidamente a cera el crecer K. Hoorl y Kennard(1970) sugieren que teles variables "no pueden soportar su potencia de predicción" y deberían sereliminadas. Con respecto a consideraciones de cómputo, la regresión por aristas es bastante eficiente ya que se puede obtener rezonablemente buenos gráficos de la traza usando solo pocos

DE DERN ARBORAL DE INGENERAL CENTRA MARRIE DE PROCESOR TRUMODO BUBLICIECA CENTRAL valores de K.La parte difícil, que es discutida en el Cap. 5 es la determinación de K,y cuen paqueña debe ser 8(K) para - justificar la eliminación de x_i. Marquandt(1974) sugiere que la eliminación de variables no es una cuestión de coeficientes iguales a cero, sinó mas bien que todas las variables debieran ser retenidas con influencia menor si 8(K) es paqueña. Esta i dea, por supuesto, ignora los motivos económicos y prácticos para la eliminación de variables.

CAPITULO No 3

CRITERIOS DE SELECCION

3.1. Funciones da criterio

Como se verá en la siguiente sección, la selección del subconjunto o subconjuntos de variables de regresión apropiados dependo del uso o que se destine la regresión. La e valuación de la información, en cada caso, será llevada a cabo mediante las funciones de criterio. Este término englobe a to dos aquellos estadígrafos que se usan para decidir si una e lección de variables es apropiada. Muchas de estas funciones son simples funciones do RSS. En esta sección se de una visión global de estos criterios.

 Media de los cuadrados de los reaiduales (Residual mean square).

2.-Cuadrado del coeficiente de correlación múltiple (Squared multiple correlation coeficient),

$$R_p^2 = 1 - \frac{R \cdot S_p}{TSS} = \frac{\widetilde{\beta}_p' x_p' y - n \widetilde{y}^2}{\gamma' y - n \widetilde{y}^2}$$

donde TSS =
$$\sum_{i=1}^{m} (\gamma_i - \overline{\gamma})^2$$

3.-R² ajustado (adjusted R²),

$$\tilde{R}_{p}^{2} = 1 - \frac{m-1}{11-p} [1-R_{p}^{2}]$$

 4.-Promedio de la varianze de predicción(average prediction variance),

5.-Total estandarizado de la media da cuadrados de los errores do estimación(standarizad total mean squared error of stimation),

6.-Promedio del la media del cuadrado del error de predicción (average prediction mean squared error),

$$s_p = \frac{RMS_p}{n-p-1}.$$

7.-Summa estandarizado de los cuadrados de los residuales (standarizado residual sum of squares),

$$RSS_{p}^{*} = e_{p}^{\prime} O_{p}^{-1} e_{p}^{\prime},$$

$$donde: e_{p} = y - \tilde{y}_{p}^{\prime} \quad y \quad D_{p} = diag(I - x_{p}(x_{p}^{\prime}x_{p})^{-1}x_{p}^{\prime}).$$

8.-Suma de los cuadredos de prediccion(prediction sum of squ<u>e</u>

Apreciamos que los seis primeros criterios son funciones simples de RSSp. Se puede argumentar que ellos son equivalentes, lo cual es cierto; pero cada criterio ofrece una perspectiva de interpretación distinta, la que debe ser considerada en forma heurística, ya que no se han daterminado las propiedades exactas de estos criterios.

3.2.<u>-Usuarios de regresión</u>

Los criterios que se usen para seleccionar el sub conjunto o subconjuntos de variables de regresión apropiados dependen del uso al que se destine la regresión. Mallows(1973)
proporciona la siguiente lista de usos potenciales de la ecua
ción de regresión:

- a)Descripción pura,
- b) Predicción simple,
- o) Extrapolación,
- d) Estimación de parámetros,

- e) Control, y
- f) Construcción de modelos

Los aspectos fundamentales que se consideran para cada uso son:

- a) Minimizar RSS (suma de los cuadrados de residuos) e la vez que se mantiene la mayor centidad de variables en el modelo. Es costumbre muy generalizada en este especto usar R², el coeficiente de correlación múltiple, en vez de RSS.
- b) Minimizar MSEP (madia del cuadrado del arror de la predicción) según las peutas fijadas en la sección 1.2.Teorema 1.1, parta (4).
- c) Minimizar MSEP con garantia de ausencia de mult<u>i</u> colinectidad. Cuando hay multicolineatidad entre las variables de predicción, la predicción que corresponde a datos fuera del rengo de ajuste suele ser ineficaz.
- d) Estimadores insesgados y de mínima varianza.Sin embargo, en algunos casos se recomianda el uso de estadísticos sesgados, siempre y cuando, ellos majoren la predicción, en especial si as una extrapolación lo que se busos.
- e) Al igual que en d) se busca estimadores insesgados de mínima varianza. En gran número de casos, hay que cons<u>i</u> derar el costo experimental.
- f) El objetivo consiste en determinar la forma o modelo de releción entre variables. Se busca la evaluación rápida de diversas ecuaciones de regresión para una elección posterior.

3.3 Evaluación de subconjuntos de regresión

Esta sección será dividida en cuatro partes. Las tres primeras son enfoques complementarios de la evaluación de los subconjuntos de regresión y la cuarta parte resume las recomendaciones de las tres primeras.

3,3,1. Aspectos básicos

En el capítulo No 1, Sección 2, Teorema 1.1, he - mos visto que si $V(\hat{\beta}_{\Gamma}) - \beta_{\Gamma} \hat{\beta}_{\Gamma}$ es semidefinida positiva, es posible estimar parámetros y respuestas con menor variabilidad utilizando el subconjunto de regresión. Escribamos $V(\hat{\beta}_{\Gamma}) = B_{\Gamma\Gamma}^{-2}$ donde $B_{\Gamma\Gamma}$ es la submatriz adecuada de $B = (X'X)^{-1}$ Una forma más operativa de verificar si $V(\hat{\beta}_{\Gamma}) - \beta_{\Gamma} \hat{\beta}_{\Gamma}$ es semidefinida positiva (sdfp) nos la da el siguiente Lema y Teorema.

Lema 3.1. See M una matriz $k \times k$ definida positiva, sean x e y dos vectores cualesquiera de k componette cada uno, entonces

Frueba. Sea p une matriz ortogonal que diagonali-

za a M, luogo
$$P'MP = D = \begin{bmatrix} d_1^2 & & & \\ & d_2^2 & O \\ & & & C \end{bmatrix},$$

con d2>0, ya que M es definida positiva.

Sea u = P'x, v = P'y ; por lo tento

x'Mxy'My - x'Myy'Mx = u'P'MPuv'P'MPv - u'P'MPvv'P'MPu

=
$$u'Duv'Dv - u'Dvv'Du$$

= $(\sum v_i^2 d_i^2)(\sum v_i^2 d_i^2) - (\sum d_i^2 v_i u_i)^2$

utilizando la desigualdad de Cauchy schwartz se tiena

Teorema 3.1. $V(\widehat{\beta}_r) - \beta_r \beta_r'$ as semidefinide positive si y solo si :

$$\frac{\beta_r' \, \beta_{rr}' \, \beta_r}{\sigma^2} \le 1 \tag{3.1}$$

<u>Fruebe</u>. Condición suficiente: sea x un vector cu<u>a</u>l quiera de r componentes. Definamos $Z = B_{rr}x$. Luago;

$$x'(v(\hat{\beta}_{r}) - \beta_{r}\beta_{r})x = z'a_{rr}^{-1}zc^{2} - z'a_{rr}^{-1}\beta_{r}\beta_{rr}^{-1}z$$
,

como 3, es definida positiva, por el lema antarior se tiene

$$x'(v(\hat{\beta}_r) - \beta_r \beta_r') \times \ge z' \beta_{rr}^{-1} z \delta^2 - z' \beta_{rr}^{-1} z \beta_r' \beta_{rr}^{-1} \beta_r$$

$$= z' \beta_{rr}^{-1} z \delta^2 (1 - \frac{\beta_r' \beta_{rr}^{-1} \beta_r}{\delta^2})$$

utilizando la hipótesis y el mecho de que 8^{-1}_{rr} es definidc positiva

$$\times'(V(\hat{\beta}_r) - \beta_r \beta_r') \times \geq Z' B_{rr}^{-1} Z S^2 \geq 0$$
.

Condición necesaria: Supongamos que $V(\hat{\beta}_r) - \beta_r \beta_r^r$ es semidefinida positiva, luego

$$\beta_r^2 \beta_r^{-1} (\beta_r r \delta^2 - \beta_r \beta_r^4) \beta_r^{-1} \beta_r \geqslant 0 ,$$
por 10 tento,
$$\beta_r^2 \beta_r^{-1} \beta_r \delta^2 - (\beta_r^2 \beta_r^{-1} \beta_r)^2 \geqslant 0 , \text{ de donde } \frac{\beta_r' \beta_r^{-1} \beta_r}{\delta^2} \le 1 .$$

For su puesto (3.1) no se puede evaluar ya que se desconoce β y 6^2 ; pero (3.1) se puede aproximar en términos del estadístico F(r,n-t-1) asociado a la hipótesiss β_Y = 0, ver Draper(1968). Se espara que

$$F = \frac{\hat{\beta}_r' = -1 \hat{\beta}_r'}{\hat{\beta}_r^2} \le \frac{1}{r} . \qquad (3.2)$$

luego, en la suposición de que el modelo es lineal y que se cum ple (3.2), es rezonable eliminar a x_r y obtener mejores estimados de los coeficientes β_p y de la respuesta y. Además, como es tos resultados son válidos para cualquier x de ingreso, estos resultados se pueden usar para extrapolación, teniendo en cuem ta las hipótesis del Capítulo 1 para el modelo.

Un criterio menos restrictivo que (3.2) es el que se sugiera en Johnston(1972) de eliminar a las variables cuyo estadístico t asociado es menor que 1. Esta condición conviene más para la predicción simple.

La condición (3.2) parece adecuada para la extrapolación moderada, pero restrictiva para la predicción. Una condición más razonable para la predicción es que $VP(\hat{y}_i)$ - $MSEP(\hat{y}_{pi})$ sea no negativa, promediando para los datos en uso se tiene:

(1)
$$\frac{1}{n}\sum_{i=1}^{n} VP(\hat{y}_i) = \frac{\sigma^2}{n}(n+t+1)$$

(2)
$$\frac{1}{n} \sum_{i=1}^{n} (VP(\hat{y}_i) - MSEP(\hat{y}_{p_i})) = \frac{rc^2}{n} (1 - \beta_r' \beta_r^{-1} \beta_r / rc^2) \ge 0$$

Demostración, si se toma $X = (X_p, X_r)$, se tiene

$$(x'x) = \begin{bmatrix} x_p^* x_p & x_p^* x_r \\ x_r^* x_p & x_r^* x_r \end{bmatrix} \equiv C = \begin{bmatrix} C_{pp} & C_{pr} \\ C_{rp} & C_{rr} \end{bmatrix}$$
$$(x'x)^{-1} \equiv B = \begin{bmatrix} B_{pp} & B_{pr} \\ B_{pr} & B_{rr} \end{bmatrix}$$

Parte (1).

$$\frac{1}{n} \sum_{i=1}^{n} VP(\hat{y}_{i}) = \frac{1}{n} \sum_{i=1}^{n} \sigma^{2} (1 + x_{i}^{2}(x^{2}x^{2})^{-1} x_{i}^{2})$$

$$= \frac{c^{2}}{n} (n + trz(x^{2}(x^{2}x^{2})^{-1}))$$

$$= \frac{c^{2}}{n} (n + trz(x^{2}(x^{2}x^{2})^{-1}))$$

$$= \frac{c^{2}}{n} (n + trz(x^{2}(x^{2})^{-1}))$$

$$= \frac{c^{2}}{n} (n + trz(x^{2}(x^{2})^{-1}))$$

Parte (2). Aplicando (1.15.e) se consigue

$$\sum_{i=1}^{n} (VP(\hat{y}_{i}) - MSEP(\hat{y}_{i})) = \sum_{i=1}^{n} (G^{2}(x_{i}^{2}\theta_{pr}\theta_{rr}^{-1}\theta_{rp}x_{pi} + x_{ir}\theta_{rr}x_{ri}) + 2x_{ip}^{2}\theta_{pr}x_{ri}) - \sum_{i=1}^{n} (x_{ip}A\theta_{r} - x_{ir}^{2}\theta_{r})^{2}$$

$$= W - \emptyset.$$

Como BprBrrBrp = θ_{pp} - C_{pp}^{-1} , deserrollando por separado estas sumas se tiene.

$$W = d^{2} \operatorname{trz} \left(X \begin{pmatrix} B_{pp} - C_{pp}^{-1} & B_{pr} \\ B_{rp} & B_{rr} \end{pmatrix} X' \right)$$

$$= \delta^{2} \operatorname{trz} \left(X'X \begin{pmatrix} \delta_{pp} - C_{pp}^{-1} & \delta_{pr} \\ \delta_{rp} & \delta_{rr} \end{pmatrix} \right)$$

$$= \delta^{2} \operatorname{trz} \left(X'XS - X'X \begin{pmatrix} C_{pp}^{-1} & 0 \\ 0 & 0 \end{pmatrix} \right)$$

como C_{pp} = X¹_pX_p, se tiene

$$W = \sigma^{2}((t+1) - trz(I_{p}))$$

= $\sigma^{2}(t+1-p) = \sigma^{2}r$.

Calculemos ahora Ø.

$$\beta = \sum_{i=1}^{n} (x_{i} + x_{i} + x_$$

como $A = -B_{pr}B_{rr}^{-1}$ se tiene

$$\beta = \text{trz} \left(x \begin{pmatrix} 3pr^{3}r^{1}\beta_{r}\beta_{r}^{i}\beta_{r}^{-1}\beta_{r}p & 3pr^{3}r^{1}\beta_{r}\beta_{r}^{i} \\ \beta_{r}\beta_{r}^{i}\beta_{r}^{-1}\beta_{r}p & \beta_{r}\beta_{r}^{i} \end{pmatrix} x^{i} \right)$$

$$= \text{trz} \left(x \begin{pmatrix} 3pr^{3}r^{1}\beta_{r}\beta_{r}^{i}\beta_{r}^{-1}\beta_{r}p & 0 \\ \beta_{r}\beta_{r}^{i} \end{pmatrix} \right)$$

$$= \text{trz} \begin{pmatrix} Cpp^{3}pr^{3}r^{1}\beta_{r}\beta_{r}^{i}\beta_{r}^{-1}\beta_{r}p + Cpr\beta_{r}\beta_{r}^{i}\beta_{r}^{-1}\beta_{r}p \\ 0 & Crp^{3}pr^{3}r^{1}\beta_{r}\beta_{r}^{i} + Crr\beta_{r}\beta_{r}^{i} \end{pmatrix}$$

$$= \text{trz} \begin{pmatrix} Cpp^{3}pr^{3}r^{1}\beta_{r}\beta_{r}^{i}\beta_{r}^{i}\beta_{r}^{-1}\beta_{r}p + Cpr\beta_{r}\beta_{r}^{i}\beta_{r}^{-1}\beta_{r}p \\ 0 & Crp^{3}pr^{3}r^{1}\beta_{r}\beta_{r}^{i} + Crr\beta_{r}\beta_{r}^{i} \end{pmatrix}$$

como CppSpr = -CprArr y CrpSpr = Ir - CrrSrr, se tiene

$$\beta = \operatorname{trz} \begin{pmatrix} 0 & 0 \\ 0 & \beta_{rr}^{-1} \beta_{r} \beta_{r}^{-1} \end{pmatrix}$$
$$= \operatorname{trz} (\beta_{r}^{+} \beta_{rr}^{-1} \beta_{r})$$
$$= \beta_{r}^{+} \beta_{rr}^{-1} \beta_{rr}.$$

finalmente se consigue

$$\frac{1}{n} \sum_{i=1}^{n} (\mathsf{VP}(\hat{y}_i) - \mathsf{MSEP}(\hat{y}_{pi})) = \frac{ne^2}{n} (1 - \frac{\beta_r^* B_r^{-1} \beta_r}{ne^2}).$$

Nótose que:

$$\sum_{i=1}^{\infty} (x_{ip}^{i} \wedge \beta_{r} - x_{ir}^{i} \beta_{r})^{2} = \beta r \beta_{rr}^{-1} \beta_{r}. \qquad (3.3.a)$$

Tomondo valores estimados a (3.3) se tiene:

$$F = \frac{\widehat{\beta}_r' \widehat{\beta}_{rr}^{-1} \widehat{\beta}_r}{F\widehat{\delta}_r^2} \le 1. \tag{3.4}$$

Por lo tento, se puede resumir que si se cumple (3.2) yel objetivo es la extrapolación o la estimación de los coeficientes β , se puede eliminar X_{r} . Por otra parte, si se cumple (3.4)y el objetivo es la predicción, se elimina X_{r} .

Un indicador del cambio de variabilidad de los parámetros estimados debido a la eliminación de variables es la ganancia relativa por predicción (Relative gain for prediction RGP)

$$RGP(\tilde{y}) = \frac{VP(\hat{y}) - MSEP(\tilde{y}_p)}{VP(\hat{y})}$$
(3.5)

que,promediando según los detos en uso, ver el Teorema 3.2, se puede escribir

$$RSP = \frac{r(1 - \beta' B^{-1} \beta_{\Gamma} / r S^{2})}{n + t + 1}$$
 (3.8)

Por lo tanto si se tione como objetivo la predicción de y,y el estimado de (3.6) RGP es \geq 0,ó en casos convenientes no es muy inferior a cero se puede eliminar X_{r} .

Resultados idénticos se obtienen para la respuesta media (estimación de y) con el único cambio del denominador de (3.6) por t + 1. Esta diferencia refleja que la variabilidad inherente del sistema supera a la variabilidad debida a la estimación de los parámetros.

3.3.2.Interpretación de Cp

Con el objeto de ilustrar los conceptos de la sección 3.2 utilizamos el estadístico C_p, ver Mallows (1973).

Como se mencionó al final de la sección 1, otros estadísti cos también puedan ilustrar la sección 3.2 en forma equiva lente a C_p . Se define a C_p como el estimador de T_p = total eg tandarizado de la media de los cuadrados de los errores de es timación de los detos X_p en uso (Standarized total meen squere error of estimation for the current data Xp)

$$T_{p} = \frac{1}{6^{2}} \sum_{i=1}^{n} MSE[\tilde{\vec{y}}_{i}]. \text{ vor (1.13.6)}$$

Teorema 3.3.

$$T_p = \frac{E(RSS_p)}{6^2} + 2p - n$$
 (3.7)

Teniendo en cuenta (1.13.g),(1.13.e) y(3.3.a) se consigue

$$T_{p} = \frac{1}{6^{2}} \sum_{i=1}^{n} (e^{2}x_{i,p}^{2} (x_{i,p}^{2} x_{i,p}^{2})^{-1} \times_{pi} + (x_{i,p}^{2} A\beta_{r} - x_{i,p}^{2} A\beta_{r}^{2} - x_{i,p}^{2})^{2}$$

$$= \sum_{i=1}^{n} [x_{i,p}^{2} (x_{i,p}^{2} x_{i,p}^{2})^{-1} \times_{pi}] + \frac{\beta_{r}^{2} \beta_{rr}^{-1} \beta_{r}}{6^{2}}$$

como:

$$B_{rr}^{-1} = C_{rr}^{-1} - C_{rp}^{-1}$$
, se tiene que

$$T_{p} = Trz(X_{p} \cdot (X_{p}^{i} \cdot X_{p}^{j})X_{p}^{i}) + \frac{\beta_{r}^{i}(C_{rr} - C_{rp}C_{pp}^{i}C_{pr})\beta_{r}}{6^{2}}$$

$$= Trz(X_{p}^{i}X_{p}(X_{p}^{i}X_{p}^{j})^{-1}) + \frac{\beta_{r}^{i}X_{r}^{i}(I - X_{p}(X_{p}^{i}X_{p}^{j})^{-1}X_{p}^{i})X_{p}^{i}\beta_{r}}{6^{2}}$$

$$+ \{n-p\} \frac{6^{2}}{6^{2}} - (n-p)$$

resordendo (1.9.b), se consigue
$$\stackrel{\mathsf{T}_p}{\mathsf{p}} = \mathsf{Trz} \left(\begin{array}{c} \mathsf{I}_p \end{array} \right) + \frac{(\mathsf{n} - \mathsf{p}) \, \mathsf{E} \left(\mathsf{RMS}_p \right)}{\mathsf{G}^2} - (\mathsf{n} - \mathsf{p}) \\ = \mathsf{p} + \frac{\mathsf{E} \left(\mathsf{RSS}_p \right)}{\mathsf{G}^2} + \mathsf{p} - \mathsf{n}$$

$$= \frac{E(RSS_p)}{6^2} + 2p - n$$

(Cuando se considere la predicción, se usa MSEP y el total ante rior suma $T_p + n$), luego C_p esta dado por

$$c_{p} = \frac{1}{63} + 2p - n. \tag{3.8}$$

Mallows(1973) proporciona indicaciones para la calibración de C_p , y recomienda que se minimice a C_p y que la condición $C_p \approx p$ indica sesgo pequeño. Como

se tiene que

$$F-1 = \frac{1}{r} \left(\frac{RSS_p}{\widehat{S}^2} - \frac{RSS}{\widehat{S}^2} - r \right)$$

$$= \frac{1}{n} \left(\frac{RSS_p}{\widehat{S}^2} - n + t + 1 - r \right)$$

$$= \frac{1}{r} \left(\frac{RSS_p}{\widehat{S}^2} - n + p \right)$$

$$= \frac{1}{r} \left(\frac{C_p - P}{\widehat{S}^2} \right)$$

por lo tento

$$C_{p}-p = r(F-1)$$
 (3.9)

Teniendo en cuenta (3.8) y (3.9) las condicènes (3.2) y (3.4) son equivalentes a :

$$p-r=2p-t-1 \le C_p \le 2p-t,$$
 (3.10)

$$p + r = 2p + t - 1 \leq C_p \leq p$$
. (3.11)

Cuando el objetivo es la extrapolación o la estimamción de parámetros y se cumpla(3.10), se recomiende eliminar

X. Para el caso de predicción se elimina X. si se oumpla

(3.11), La expresión (3.11) es consistente con la recomendación

de que C_p see pequeño o cercano a p.

La relación de C_p con RGP la obtenemos por medio de (3.5)

$$\widehat{RSP} = \frac{P - C_{P}}{D + t + 1}$$
 (3.12)

Fara el caso en que un C_p sea mayor que p, al eliminar X_p se produce una pérdida de exactitud. En este punto, si el objetivo es la predicción, se recomienda buscar la combinación de p elementos que minimize C_p y se elimina su X_p asociado. Fara el caso de estimación de respuesta media, la expresión ASP se obtiene al eliminar n en el denominador de (3.12).

En forma análoga se deducen otras aplicaciones de \mathbb{C}_p equivalentes a las usadas con otros estadísticos for e jamplo, la condición $f\leqslant 2$ es equivalente a $\mathbb{C}_p\leqslant t+1$.

3,3.3. Otras funciones de criterio.

gráfico, son analizados muy frecuentemente.La selección de pose hace considerando tres aspectos[i] mínimo RMSp, (ii) el valor de potal que RMSp = RMS o (iii) el valor de potal que el lugar galmétrico del mínimo RMSp crece brúscamente cuando podecrece.Do (3.8) se tiene

$$C_p = (n-p)(\frac{RNS_p}{\hat{\xi}^2} - 1) + p$$
, RMS = $\hat{\xi}^2$ (3.13)

Esta fórmula indica una relación directa entre C_p y RMS_{p.} a relación (3.10), que indica eliminación de X_p cuando el objetivo es extrapolación y estimación de parámetros, se expresa camo

$$\frac{n-t-1}{n-p} \quad \text{FMS} \leqslant \text{FNS}_p \leqslant \frac{n-t}{n-p} \quad \text{FMS} \tag{3.14}$$

Les condiciones (ii) y (iii) se basan en que -----

RMS $_p$ /RMS \approx 1. En este caso $C_p \approx P_p$ pero hay que tener cuidado en que C_p magnifica la diferencia del cociente con la unidad por el factor (n-p) lo cual hace muy sensible a C_p a los cambios de RMS $_p$ /RMS. Las condiciones (ii) y (iii) se recomiendan para la eliminación de X_p cuando el objetivo es la predicción, tenien do siempro en cuenta la discusión de (3.12).

 $\frac{R^2}{P}$; esta critario probablemente es uno de los más utilizados y su gráfico es empleado para la determinación de p,a medida que p decrece el lugar geométrico del máximo $\frac{R^2}{P}$ se mantieno constante hasta un p en el que decrece brúscamente, esto indica el p adecuado. La relación entre $\frac{R^2}{P}$ y $\frac{C}{P}$ es

$$C_p = \frac{(n-t-1)}{1-R^2} + 2p-n$$
, (3.15)

esto indica que pequeñas veriaciones de \overline{R}_p^2 son más sensibles por C_p debido al factor de amplificación (n-t-1). Por ello, el criterio \overline{R}_p^2 indica la eliminación de más variables que C_p . Se critica al estadístico \overline{R}_p^2 porque hay casos en los que se elimina variables importantes. También se presentan dificulta des para extraer conclusionos a partir de los gráficos y, como alternativa, se sugiero a \overline{R}_p^2 ajustado: \overline{R}_p^2 . \overline{R}_p^2 sugiero elegir el para el que \overline{R}_p^2 es máximo. Su relación con RMS $_p$ es

$$\vec{R}_p^2 = 1 - (n-1) \frac{RMS}{TSS}, (TSS = Y'Y - n\overline{Y})$$
 (3,16)

y es exéctamente equivalente a buscar el mínimo de RMSp.

 $\frac{J_p}{p}$; se obtione al computer el total de la varianza de predicción en el conjunto de datos X_p en uso y de la estimación de 6^2 por RMS. Este estadístico tiene la desventaja de ignorar al sesgo (ver 1.14.c).

 $\frac{S_p}{p}$; tionon una definición silmilar a C_p - S_p se cri-

gine a partir del promedio de la media del cuadrado del error de predicción de los datos X en uso, donde X y Y se distribuyen según una normal multivariente, ver[1.14.d] en este caso se tiene que

$$E(MSEP) = \frac{(n+1)(n-1)}{n} \frac{E(RMSp)}{n-p-1}.$$
 (3.17)

El estadístico S_p se obtiene de eliminar el factor<u>(n+1)(n-1)</u>

y estimar E(RMS_p) por RMS_p Aceptando la hipótesia de normalidad multivariante, en predicción as sugiero utilizar la ecuación que tenga mínimo S_p y, como el promedio se toma sobre todos los datos, se puede usar el conjunto con S_p mínima para extrapola - ción moderada.

Para el problema de predicción y control, Lindley

[1968] desarrolla una aproximación bayesiana al problema de se
lección de variables, incluyendo en su análisis el costo de obser

vación de les variables X. En este estudio se recomienda minimi

zer la cantidad

$$\frac{\hat{\beta}_r' B_{rr}^{-1} \hat{\beta}_r}{D} + U_p, \qquad (3.18)$$

donde U es el costo de observación.Anteriormente yo se ha re -comendado el subconjunto que minimiza el primer sumando.

Pora el problema de control de la salida el rededor de un punto y_0 , Lindley recomisada seleccionar el subconjunto de p variables que minimiza

$$\frac{\hat{\beta}_{r}' \, \beta_{rr}^{-1} \, \hat{\beta}_{r}}{n} + u_{p} + 6^{2} \left(\frac{r}{n} + \frac{y_{o}^{2}}{6^{2} + \hat{\beta}_{p}'(x_{o}' x_{p}) \hat{\beta}_{p}} \right). \quad (3.19)$$

3.3.4.-Recpmendaciones para la selección de variables

A continuación se da un resumen de recomendaciones para la selección de variables, según los usos de la ecuación

- quo so mencianan en la segunda Sección de este Capítulo.
- a)Para la descripción pura, se recomiende el subconjunto que minimiza RSS o que maximiza R², eliminándose popes varia blos.
- b)Pera la predicción se recomienda el subconjunto que cumple con(3.4) o (3.11),o que hace cero (3.6),o que hace cero (3.12),o RNS_p = RNS,o que min.C_p,o máx.R_p²,o min.S_p,o $max \frac{R^2}{p}$.

Para la estimación de respuesta se recomienda hacer cero la modificación de (3.6)ó (3.12),o máx. ${
m R}_{\rm p}^2$,o máx. ${
m R}_{\rm p}^2$.

- c)Para la extrapolación se escoje el subconjunto que cumple (3.2),**0** (3.10) o (3.14),o que min.5_o.
- d) Fara la estimación de parámetros se recomienda que el subconjunto de veriables satisfaga (3.2),o (3.10),o (3.14) o que máx. R_p^2 o \tilde{R}_p^2 .
- e)Para el caso de control se recomienda satisfacer(3,18) o (3,19).
- f)En modeleje se recomiande eyudarse con esesoría técnica y con algoritmos de cómputo eficietes y fijerse alguno de los usos anteriores pare la ecuación de regresión.

3.4. Selección de variables según los métodos de regresión por pasos

En esta sección veremos los conceptos de selección para los métodos de selección hadía adelanta(FS) y eliminación hacia atras(BE). Según se vió en el Capítulo No 2, los cálculos son secuenciales. Un punto bastante delicado es la determinación del criterio de finalización, según el uso de la ecua ción de regresión mediante determinación de $F_{in}=F(\times,1,n-p-1)$ é

 $F_{\rm out}=F(\propto,1,n-p)$. En lineas generales, en SE si $F_{\rm in}$ es pequeño o en FS si $F_{\rm out}$ es pequeño los modelos incluyen numerosas variables; lo contrario sucede si $F_{\rm in}$ ó $F_{\rm out}$ es grande. Se puede recomendar un proceso consistente en procesar SE o FS en su tota lidad, obtener un conjunto de cada tamaño y seleccionar al final.

En Hocking(1975) se menciona que en un estudio de - simulación, Bendel(1974), compara ocho reglas distintas de finalización para FS. Su estudio incluye a C_p , S_p , F univariada, F.falta de ajuste (L_f)

$$L_{f} = \frac{RSS_{p} - RSS}{\hat{\varsigma}^{2}(t+1-p)},$$

 \widetilde{H}_p^2 y otros;y se recomiendo el uso de F univariado o la prueba $S_p = S_p$ con nivel $0.1 \leqslant \alpha \leqslant 0.4$ para pocos grados de libertad Para 40 ó mas grados de libertad, C_p y S_p son los mejores,y F univariada es bastante bueno para .10 $\leqslant \alpha \leqslant$.25. Se sugisre que el mejor test es F univariada con $\alpha = 0.15$.

Se consigue reducir p si se utilizan métodos de selección mas eficaces. Por ejemplo, si se evalúan todos o al menos los mejores conjuntos.

Para la regresión polinomial se recomienda tener precaución con la regla común de finalización, cuando la siguien ta potencia no mejora la ecuación.

Finelmente, se puede obtener información sobre eva luación y verificación de una experiencia en Hocking(1976).

CAPITULO No 4

MULTICOLINEALIDAD

4.1. Aspectos generales

Fera resolver ol modelo de regresión lineal $Y = X \beta + \mathcal{C} \tag{7.1}$

Con t + 1 variables x_i independientes y una variable y dependiente por mínimos cuadrados, se requiere que el rango de la matriz de datos X sea t + 1, o equivalente, $|X'X| \neq 0$ (en este caso X'X es definida positiva) y la estimación de β es $\hat{\beta} = (X'X)^T X'Y$. Si el modolo es correcto (en el sentido de incluir a todas las variables independientes de importancia) se tiene

 $E(\hat{\beta}) = \beta$, $V(\hat{\beta}) = 6^2(X'X)^{-1}y$ para la predicción simple $V(\hat{y}(x)) = 6^2(1 + x'(X'X)^{-1}x)$. Esto nos indica que la varia bilidad de $\hat{\beta}$ y de \hat{y} está dada por la magnitud de las componentes de $(X'X)^{-1}$, en sentido más amplio por una norma de $(X'X)^{-1}$. For lo tento, es de importancia que las componentes de $(X'X)^{-1}$ tengen una cota rezonablemente baja. Esto se consigue cuando (X'X) no es ceroana a cero.

Veamos, entonces con más detenimiento el det.de X'X Si |X'X| = 0, se tiene que una o más variables independientes se pueden deducir de las otras mediante combinaciones lines - les y que no se puede estimer β de (4.1). Vemos, pues, que si algunas variables independientes se pueden deducir linealmente, de las otras, esto imposibilita la deducción de los parámetros β .

Si $|X'X| \approx 0$ se tiene que las variables independientes tienen una relación cosi lineal y que V $(\hat{\beta})$ y $V(\hat{y})$ pueden ser muy grandes.

Si $\{X'X\} \gg 0$, es el caso más Favoreble. Esta condición se logra cuando $X_i^*X_j=0$, e sea cuando las columnas de $X_i^*X_j=0$ aon ortogonales.

A los conceptos anteriores se los estudia bajo el nombre de multicolinealidad.

Definición 4.1 : Se dice que les variables independientes X de $\{7.1\}$ son completemente multicolineales si $\|X^*X\| = 0$, y que hay falta completa de multicolinealidad cuando do $X_1^*X_2^* = 0$. Entre estos dos casos extremos, tenemos diversos grados de multicolinealidad, la cual no se refiere únicamente a relaciones de tipo lineal entre las variables independientes.

En general, la existencia de multicolinealidad en la matriz X devienc en: (1) la estimación inexacta de los cogficientes de regresión, tanto por la varianza grande como por la inestabilidad numérica de la solución: (2) la especificaci en incierta del modelo con respecto a la inclusión de variables; y (3) dificultad en determinar en que medida las variables independientes influyen en la variable dependiente y.Del primer punto hemos tenido una idea clara en esta sección y para los otros dos veamos primero los conceptos de regresión ortogonal y de regresión auxiliar.

4.2. Regresión ortogonal

Si las columnas de X son ortogonales dos a dos, se tiene

$$X_1' X_2 = 0$$
 (4.2)

para cualquier partición X_1X_2 de X. Sin pérdide de generalidad se puede asumir que $X = (X_1, X_2)$. Luego $\hat{\beta}$ esta dedo por

$$\hat{\beta} = \begin{pmatrix} \hat{A} \\ \hat{\beta}_{2} \end{pmatrix} = \begin{bmatrix} (x_{1} \cdot x_{1})^{-1} x_{1}^{2} \\ (x_{2}^{2} \cdot x_{2})^{-1} x_{2}^{2} \end{bmatrix}$$

$$(4.3)$$

Se eprecia una independencia completa de $\hat{\beta}_1$ con $\hat{\beta}_2$ y que $\hat{\beta}_1$ se puede obtener regresionando Y con X₁ y que al agregar X₂ y regrasionar Y con X = (X₁, X₂), $\hat{\beta}_1$ no cambia, y lo mismo sucede para $\hat{\beta}_2$

Cualquier discrepancia de la ortogonalidad dos a dos de las columnas de X, indica la presencia de multicolina<u>m</u> lidad, la cual se acentúa a medida que |X'X| tiende a cero.

Téngase presente que el método de mínimos cuadrados no es invalidado por la presencia de multicolinealidad en X. Este método solo fella cuando el rango de x es menor que t + 1; lo que sucede es que ______los datos no permiten distinguir su influencia sobre la variable dependiente y .

4.3. Efectos de la multicolinealidad en la especificación del modelo

Supongamos que el modelo de regrasión (4.1) de t + 1 variables es el adecuado, pero que X tiene un alto grado de multicolinación por lo que los estimados obtanidos son no eignificativos y se duda de la especificación original. Una alternativa para eliminar la multicolinealidad es la eliminación de les variables que la acentúan. Usemos el esquema de Huang(1970) para ver el efecto de esta restricción.

Supongamos que (4.1) representa el modelo adecuado pero que en realidad la ecuación se ajusta según

$$Y = \bar{X} \bar{\theta} + \bar{\mathcal{C}}, \qquad (4.4)$$

donde \overline{X} se una submatriz de \overline{X} , \overline{X} tiene m columnes, m < t+1. Por mínimos cuadrados, se tiene que

$$\widehat{\beta} = (\overline{X}'\overline{X})^{-1}\overline{X}'Y = (\overline{X}'\overline{X})^{-1}\overline{X}'(X\beta + Q)$$
 (4.5)

luego

$$\mathsf{E}(\hat{\boldsymbol{\beta}}) = \mathsf{Q}\boldsymbol{\beta}, \tag{4.6}$$

donde

$$Q = (\tilde{X}'\tilde{X})^{-1}\tilde{X}'X = ((\tilde{X}'\tilde{X})^{-1}\tilde{X}'X_1, \ldots, (\tilde{X}'\tilde{X})^{-1}\tilde{X}'X_{t+1}).$$
 Apreciamos que una columna j cualquiera de Q está formada por coeficientes de la regresión de x_j como variable dependiente con \tilde{X} como variables independientes. Por lo tento, Q_j indica - como influye \tilde{X} en la determinación de x_j , según el modelo

$$x_{i} = q_{1j} \bar{x}_{1} + q_{2j} \bar{x}_{2} + ... + q_{mj} \bar{x}_{m} + U_{j}$$
 (4.8)

A esta relación, que asocia la variable "verdadera" x $_j$ con las variables incluidas \tilde{x}_4 , \tilde{x}_2 , ..., \tilde{x}_m , se le llama regresión aux \tilde{x}_1 liar. For lo tanto, si se tiene un concepto apriori de lo que es \tilde{q} , según (4.6) se pueden deducir ciertas características de $\hat{\beta}$.

Fara ilustrar estas ideas, veamos un ejemplo con - creto. Supongamos que \widetilde{X} se obtiene a partir de X por la elimina ción de una columna, por ejemplo X_{t+1} . Entonces,

$$Q = (I_t, q), q' = (q_1, q_2, ..., q_t)$$
 (4.9)

según (4.6) so tiene:

$$E[\hat{\beta}_{1}] = \beta_{1} + q_{1}\beta_{t+1},$$

$$E[\hat{\beta}_{2}] = \beta_{2} + q_{2}\beta_{t+1},$$

$$E[\hat{\beta}_{1}] = \beta_{t} + q_{t}\beta_{t+1}.$$
(4-10)

El producto $q_i \beta_{t+1}$ de el sesgo de $\hat{\beta}_i$; q_i indica la influencia de \bar{X}_i sobre X_{t+1} , y β_{t+1} es la inf. de X_{t+1} sobre Y, por ejemplo en la determinación de una producción β_i según las variables independientes : labor L y capital C, se espera la inclusión de la variable dirección gerancial β_i . Luego, si regresionemos β_i con L y C únicamente, se espera un sesgo en β_i y β_i cuyas direcciones, positiva o negativa, se deducen teniendo en cuenta que la influzencia de β_i β_{k+1} sobre

p es positiva y las influencias de L y C $\{q_{_L},q_{_L}\}$ sobre D son positivas. For lo tanto, en la regresión de C y L sobre P se espora seegos positivos en $\beta_{_C}$ y $\beta_{_L}$.

4.4. Detección de multicolinealidad *

En la mayoria los modelos, los datos presentan cierto grado de multicolinealidad, y el objetivo general no es eliminarla sino reducirla. Una forma de identificar la dependencia entre dos variables es a través de sus coeficientes de correlación, los cuales estan dados por los componentes de la matriz M = X'X, y una inspección de ella nos permite ver la correlación de las variables dos a dos. El caso de felta completa de multicolinealidad se de para datos ortogonales y se tiene M = I (por lo tanto |M| = 1); y el otro extremo se de cuando |M| = 0. Para los casos intermedios, en Huang(1970) se sugiere que una condición "tolerable" de multicolinealidad - entre xi y xi se de cuando

$$|M_{ij}| < R = \sqrt{R^2}$$
 (4.11)

Pero , es el usuario, el que en definitiva acota M_{ij} 1. Nótese que este proceso sirve para analizar grupos de dos variebles eolemente, y que su aplicación por transitividad, para mas 'de dos variables, ya no es confiable, en este caso se recomienda el coeficiente de correlación múltiple R_i^2 o el esta tadístico F_i de cada variable \mathbf{x}_i con las otras t restantes Un elevado valor de R_i^2 o de F_i indica un alto grado de relación de \mathbf{x}_i con las otras variables.

^{*} Par facilidad, en esta sección se trabajará con los datos tomados como desviaciones de su valor medio y divididos por su desviación estandard. A esta proceso se la llamaratandarización de los datos.

Algunos conceptos muy importantes se den en Gunst (1977), al detectar la multicolinealidad con ayuda de valores y vectores propios. Sean λ_i y Z_i los valores y vectores unitarios propios de M = X'X. Las λ_i cumplen $\Sigma \lambda_i$ = t + 1 Y, para el caso de ortogonalidad, se tiene λ_i = 1 para todo i. Formamos la matriz $Z = (Z_1, Z_2, \dots, Z_{t+1})$. Luego, a partir de $V(\widehat{\beta}) = \delta^2(x'x)^{-1}$, se obtiene

$$V(\hat{\beta}_{j}) = \hat{\epsilon}^{2} \left(\frac{z_{j1}^{2}}{\lambda_{1}} + \frac{z_{j2}^{2}}{\lambda_{2}} + \dots + \frac{z_{j_{t+1}}^{2}}{\lambda_{t+1}} \right), 1 \le j \le t+1,$$
(4.12)

lo cual indica la importancia de que λ_j sea grande.Suponga - mos que $\lambda_j \approx 0$, luego X'XZ $_j = \lambda_j$ Z $_j \approx 0$, se puede esparar que XZ $_i \approx 0$, o sea

$$\sum_{i=1}^{t+1} X_i Z_{ij} \approx 0. \tag{4.13}$$

En esta combinación lineal, los coeficientes Z_{ij} con mayor ve lor absoluto indican que las variables correspondientes x_i es ten influyendo más fuertemente en la multicolinealidad.

En Gunet(1977) se presenta un modelo de una varia ble y, con 24 variables x_i multicolineales, se calculan los tros valores propios mas pequeños $\lambda_i \in \lambda_2 \leqslant \lambda_3$ y sus vectores propios asociados Z_1 , Z_2 , y Z_3 (ver table 4.1). Se aprecia que

 λ_1 = 0, indica completa multicolinealidad entre les variables X_i cuyos $\{Z_{i,1}\}$ correspondientes son grandos, en este caso x_0 hasta x_{20} .

 λ_{2} = .060, indica multicolinealidad fundamentalmenta entre x_{4} , x_{5} , x_{8} , y x_{24} .

 λ_3 = .121, esta λ todavia es pequeño comparado con la unidad (en ortogonalidad λ_i = 1 para todo i) y puede indicar

multicolinealidad, pero es recomendable reducir por eliminación de variables la multicolinealidad indicada por λ_4 y λ_2 y des pués en el modelo reducido enalizar nuevamente la multicolinea lidad.

4.5. Reducción de la multicolinealidad *

Una técnica muy utilizada es la imposición de restricciones lincales en los coeficientes β, tales como la elimi nación de coaficientes y otras restricciones.

Supo**ngemo**es q**ue β cumpl**e

$$p \beta = Q \qquad (4,14)$$

donde P es una matriz (t+1) x (t+1) conocida y $\mathbb R$ un (t+1) vector conocido, luego el estimador $\widetilde{\beta}$ de β que édmpla. (4.14) y minimiza los cuadrados de los residuales es $\widetilde{\beta} = \widehat{\beta} + (x'x)^{-1}P'(P(x'x)^{-1}P')^{-1}(Q - P\widehat{\beta}), \qquad (4.15.4)$

$$\beta' = \beta + (x'x)^{-1}P'(P(x'x)^{-1}P')^{-1}(Q - P\beta),$$
 (4.15.2)

$$E[\vec{\beta}] = \beta + (x'x)^{-1}P' \cdot (P(x'x)^{-1}P')^{-1}(Q - P\beta)$$
 (4.15.b)

$$V(\tilde{\beta}) = G^{2}(X'X)^{-1} - G^{2}(X'X)^{-1}P'(F(X'X)^{-1}P')^{-1}F(X'X)^{-1}.$$
(4.15.5)

3i $P\beta = Q$, $\widetilde{\beta}$ es insesgado y de (4.5) vemos que $V(\widetilde{\beta}_i) \leq V(\widehat{\beta}_i)$, cuando la restricción es aproximadamente cierte $\stackrel{\sim}{eta}$ es sesgado y la medida de variabilidad que nos intesesa es MSE(eta)

MSE(
$$\vec{\beta}$$
) = E($\vec{\beta}$ - β) ($\vec{\beta}$ - β)

$$= 6^{2} (I - (X'X)^{-1} P'(P(X'X)^{-1} P) (X'X)^{-1}$$

$$+ (X'X)^{-1} P'(P(X'X)^{-1} P')^{-1} (Q - PB)(Q - PB)'(P(X'X)^{-1} P')^{-1} P(X'X)^{-1}.$$

le discrepancia de la restricción $F\beta = Q$ no es significat<u>i</u> ve, a menudo se tiene que MSE $(\widetilde{m{eta}}_1)$. \leqslant $V(\widehat{m{eta}}_1)$. Para ilustración consideremos el modelo

cuyos vectores de observación de x_1 y x_2 son X_1 y X_2 respectiva-mente.

* En esta sección se trabaja con datos estendarizados.

TABLA 4.1

Vectores propios da X'X correspondientes a los tres valores propios más pequeños.

,	λ ₁ = 0	λ ₂ =.060	x3=.121
Variable	z ₁	^z a	z ₃
×1	ο.	:003	0 53
× ⁵	0.	.120	.374
×3	o.	.105	-,631
×4	ο.	337 [*]	.204
× ₅	a.	.487*	.118
× _G	ο.	.0768	. 1 59
× ₇	ο.	093	291
×e	0.	.425 [*]	.321
' 's	.497	003	.091
×10	.300	.022	036
×11	. 270	001	025
×12	.349	.015	019
× ₁₃	. 325	010	067
×14	.139	.055	034
×15	. 199	.019	.074
×16	.195	003	107
×17	.139	034	.053
× _{1a}	.348	.038	. 08 0
× ₁₉	.270	038	077
×20	.270	011	004
×21	0.	.030	168
, ss	ο.	-,136	.167
× ₂₃	٥.	.01ឆ	3 59
×24.	σ.	-,631 [*]	. 15 9

Lusgo:

$$x = (x_1, x_2),$$

$$x'x = \begin{pmatrix} 1 & P \\ P & 1 \end{pmatrix}$$

$$(\hat{\beta}_1) \quad \{^2 \quad \{^1 \quad P\}\}$$

У

$$V\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \frac{g^2}{1 - \rho^2} \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} ,$$

donde ρ es el coeficiente de correlación entre $x_{1,y}$ x_{2} . Asumamos la restricción β_2 = 0 cuya notación metricial es

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

$$NS4(\widetilde{\beta}) = \frac{6^2}{1-p^2} \begin{pmatrix} 1-p^2 & 0 \\ 0 & 0 \end{pmatrix} + \beta_2^2 \begin{pmatrix} p^2 & -p \\ -p & 1 \end{pmatrix}$$
 (4.17.a)

Comparemos

$$V(\hat{\beta}_1) = \frac{\varsigma^2}{1 - p^2}$$

COD

MS(
$$(\vec{\beta}) = \delta^2 + \beta_z^2 e^2$$
. (4.17.b) Fare algunes valores do β y β_z , to contided (4.17.b) es menor que (4.17.a).

La restricción lineal mas utilizada es la restricción cero, o sen que, a algunos β_i se les iguals a cero, Esta tes β_i son alegidos de acuerdo a que las columnas X_i assoladas tengan elta multicolinealidad con las otras columnas de X y poca relación con el vector y (se analizan los estadísticos F_i , R_i^2 , etc. Una forma muy apropiada de esta proceso se presento an Eunet(1977). La idea es la siguienta :

Se forma la matriz A=(Y,X). Séan $\bar{\lambda}_1$, \bar{Z}_1 los valores y vectores propios de A^*A_1 , según la discusión de la sección anterior. Si existe $\bar{\lambda}_j^* \approx 0$, esto indica multicolinamidad. Las componentes \bar{z}_{ij} del vector propio \bar{Z}_j con mayor valor indicam que las columnas A_i producen multicolinamidad. Si la columna A_1 de A_1 no está incluida entre las A_i esto nos indica que la

multicolinealidad se da principalmente entre las columnas de X, y que éstas no influyan mayormente en y, y por lo tanto se pug den climinar. Por ejamplo Gunet(1977) presenta la tabla 4.2 (página siguiente).

 $\bar{\lambda}_{\rm s}=0$ indice multicolinealidad complete entre las columnas ${\rm A_i}$ cuyas ${\rm IZ_{io}}$ correspondientes son grandes, en este caso ${\rm x_q}$ hasta ${\rm x_{20}}$ y no se influencia a y, ya que ${\rm Z_{in}}=0$.

 $\bar{\lambda}_1$ = .059 indica multicolinealida entro X_4 , X_5 , X_8 , y X_{24} y no se influencia a y (\bar{Z}_{01} = .055). Identicamente se puede inspeccionar $\bar{\lambda}_2$, el cual es relativamente mayor que $\bar{\lambda}_0$ y $\bar{\lambda}_1$.

El análisis enterior nos sugiere eliminar variables entre $x_q = x_{20}$ y entre x_4 ó x_5 ó x_8 ó x_{24} . Otros criterios esta dísticos (t,f, etc) nos ayudan a elegir. Repitiendo este proceso dos veces se obtiene la table 4.3 con menos datos.

Pare este caso $\lambda_0=0.100$ puede ser considerado cercano a cero, y las veriables que exhiben multicolinealidad - fuerte son Y, $X_5, X_8, y X_{24}$. Como en la multicolinealidad de incluye a y, se recomiende detener la eliminación. En esta forma se ha conseguido un modelo con pocas variables, en el que se reduce la multicolinealidad entre las X_i y, además, se manifiesta una relación fuerte de las X_i retenidas con Y.

TABLA No 4,2

Tres valores propios de A'A y sus correspondientes vectores propios.

praptos.			
	$\bar{\lambda}_{\sigma} = .0$	ર્ગ=.059	$\bar{\lambda}_2 = .115$
Variable	ই ০	₹1	Že
У	0,	:056	.274
×1	ο.	-,010	127
ײ	ο.	.124	.242
×э	٥.	.097	-,547
×4	ο.	3 32*	.215
× ₅	0	.475 [*]	-,050
×s	ο.	.076	.099
× ₇	Ο.	100	269
×e	O	.419 [*]	.265
×a	.497	001	.111
×10	.300	.019	051
×11	.270	004	045
×12	.348	.020	009
×13	.325	~.008	041
×14	.139	.023	024
×15	.139	.024	.097
×16	.195	010	098
×17	.139	058	-100
×19	.348	.034	.034
×19	.270	~.041	058
×so	.270	0 21	045
×s1	0.	.012	258
×ss	0.	-,125	.215
×za	0.	.001	413
×24	c.	648 [*]	.035

TABLA NO 4.3

Los valores propios de A'A con sus correspondientes vectores propios para el modelo reducido de 15 variebles.

•		
	λ ₀ = . 100	278-156
Variable	z _o	z ₁
У	152 [*]	648
×1	.042	. 17 8
× ₃	033	143
×s	a39 *	545
×e	067	.192
×7	,000	.024
×a	498 [*]	.114
×10	012	.118
×12	004	.043
×17	.036	122
×18	015	.281
××30	.070	.105
×21	.001	.273
×ss	.033	127
×s ³	.019	.179
×24	.911	.108

CAFITULO No 5

REGRESION POR ARISTAS *

(Ridge Regression)

5.1. Aspectos generales

La forma usual de estimar los perámetros β , en un proceso de ajuste de curvas, so hace por la técnica de mínimos cuadrados, la cual es buens cuando cuando X'X no discrepa considerablemente de I. En este caso, los valores propios λ_i de X'X no son cercanos a cero; pero, cuando X'X es muy distinta de I, ocurron problemas de inflación e inestabilidad general. Hoerl(1962), sugirió perturbar la disgonal de X'X con el objeto de aproximarla a I y estimar β , según

$$\hat{\beta}^* = \{x'x + \kappa I\}^{-1} \ x'Y$$

$$= WX'Y \ W = \{x'x + \kappa I\}^{-1}$$
(5.1.a)

A la estimación y el análisia construidos de acua<u>r</u> do a (5.1) se la llama "Regresión por aristas", y su relación con la estimación ordinario está dado por

$$\hat{\beta}^* = (\kappa(x^*x)^{-1} + I^{-1}) \hat{\beta} = Z \hat{\beta} ,$$

$$Z = (\kappa(x^*x)^{-1} + I)^{-1} = (x^*x + \kappa_I)^{-1} x^*x . \quad (5.1,6)$$

Fropiedades mas importantes de P* W y Z

(i) Sean $f_i(W)$ y $f_i(Z)$ los valores propios de W y Z respectiv<u>e</u> mento,, antonces

$$\xi_{i}(w) = \frac{1}{\lambda_{i} + \kappa}$$
 (5.2.a)

$$\xi(z) = \frac{\lambda_i}{\lambda_i + K} , \qquad (5.2.6)$$

* En este Capítulo se trabaja con datos estendarizados.

donde λ_i son los valores propios de X'X. Estas relaciones se obtienen a partir de (5.1.a) y (5.1.b).

(ii)
$$z = I - K(x'x + KI)^{-1} = I - KW,$$
 (5.9)

Esta resultado se obtiene tomando la forma alternativa de $Z = WX^{3}X$.

(iii) Para K \neq 0, $\widehat{oldsymbol{eta}}^{oldsymbol{\mathsf{t}}}$ esto es

$$(\hat{\beta}^*)$$
, $\hat{\beta}^* < \hat{\beta}'$ $\hat{\beta}$ (5.4)

<u>Prueba.</u> See [.] le norme Euclidea en \mathbb{R}^n , see [] le norme espectral, en $M(n \times n)$; en elgebre se establece que establementables, competibles, este quiere decir que , para toda matriz M y todo vector X se cumple $|MX| \leq ||M|| \cdot |X|$. Z es definide positive. For lo tento,

$$(\widehat{\beta}^*)^{'}\widehat{\beta}^* = (z\,\widehat{\beta})^{'}z\,\widehat{\beta} = |z\,\widehat{\beta}|^2 \leqslant |z\,||^2 |\widehat{\beta}|^2 = \xi^2_{\max}(z)\,\widehat{\beta}^{'}\widehat{\beta}.$$

Como todo $\zeta_i(z)$ es menor que la unidad, (5.4) se verifica fácilmente.

(iv) El estimador de la suma de los cuadrados de los residuales esta dado por

ครอ(
$$\hat{\beta}$$
) = (Y - \hat{x}) (Y - \hat{x}) (Y - \hat{x}) = Y'Y - ($\hat{\beta}$) X'Y + ($\hat{\beta}$) X'X $\hat{\beta}$ * - ($\hat{\beta}$) X'Y = Y'Y - ($\hat{\beta}$) X'Y - \hat{x} ($\hat{\beta}$) $\hat{\beta}$ * (5.5)

= (suma total de cuadrados - suma de cuadr<u>a</u>

dos debido a la regresión de β - cuadrado

de longitud de β por K)

Observación: Da (5.2.6) y (5.3) sa tiena,

$$Z(0) = I$$

5.2. La traza de aristas

5.2.1 <u>Definición de la traza de aristas</u>

A medida que X'X se desvía de la matriz identi - dad I, o sea, cuando se tiene valores propios pequeños, la probabilidad de que $\hat{\beta}$ esté cerca a $\hat{\beta}$ puede ser muy pequeña $\left(\vee \left(\hat{\beta} \right) = e^{2} (x'x)^{-1} \right)$.

La inspección de la correlación entre las varia — bles tomadas dos a dos (elementos de X'X) no es suficiente para determinar las relaciones existentes para más de dos facto res. Los métodos computacionales actuales no permiten un conocimiento adocuado del espacio de factores y de la sensibilidad de los resultados, en cada caso particular, con excepción de los métodos de componentes principales y regresión por aristas. Se obtiene información más tengible de la traza de las aristas esto es : calcular $\hat{\beta}^*(k)$ y $853(\hat{\beta}^*)$ para varios K.

5.2.2 Caracterización de la treza de aristas

Sea 8 cualquier estimador lineal de β , la suma del cuadrado de los residuales es

RSS(3) =
$$(Y - XB)^{1}(Y - XB)$$

= $(Y - X\hat{\beta} - (XB - X\hat{\beta}))^{2}(Y - X\hat{\beta} - (XB - X\hat{\beta}))$
= $(Y - X\hat{\beta})^{2}(Y - X\hat{\beta}) + (B - \hat{\beta})^{2}(X^{2}X(B - \hat{\beta}))$
= $(B)^{2}$

La traza de aristas ' se define como la 8 que re - suelve el problema siguiente:

Sujoto a (3-
$$\hat{\beta}$$
)'x'x(3- $\hat{\beta}$) = ϕ_{o} .

Utilizando multiplicadores de Lagrenge con multiplicador - 1,

(5.8.a) es equivalentera

Minimizar
$$f = F(B,K) = B'B + \frac{1}{K}(B - \hat{\beta})'X'X(B - \hat{\beta}) - \beta$$
, (5.8.6)

cuya solución es justamente la traza de aristas

$$B = \hat{\beta}^* = (x'x + kI)^{-1}x'Y,$$
 (5.8.0)

con K elegida de tal forma que se satisface la restricción (5.8.a)

- 5.3 <u>Proisdades del error cuadrátrico medio de la regresión</u>
 por eristes
- 5.3.1 Varianza y sesgo de un estimador por aristas.

Definemos :

$$L_{1}^{2}(K) = \begin{cases} \text{Cuadrade de } \beta \\ \text{distancia de } \beta \\ \text{e} \end{cases},$$

$$L_{1}^{2}(K) = (\beta^{*} - \beta), (\beta^{*} - \beta). \qquad (5.9.a)$$

Aplicando el operador esparanza se tiene :

$$\begin{split} & \in (L_{f}^{2}(K)) = E(I \hat{\beta}^{*} - \beta)'(\hat{\beta}^{*} - \beta)) \\ & = E(I \hat{\beta} - \beta)'Z'Z(\hat{\beta} - \beta)) + (Z\beta - \beta)'(Z\beta - \beta) \\ & = e^{Z}(Irz Z(X'X)^{-1}Z') + \beta'(Z - I)'(Z - I)\beta, \end{split}$$

utilizando (5.1.b) y (5.1.a)

$$E(L_{1}^{2}(k)) = 6^{2}(\text{Traze}(X'X + KI)^{-1} - K \text{ traze}(X'X + KI)^{-2})$$
$$+ K^{2}\beta'(X'X + KI)^{-2}\beta$$

aplicando (5.2.a)

$$E(l_{1}^{2}(k)) = 6^{2} \sum_{i} \frac{\lambda_{i}}{(\lambda_{i} + K)!} + K^{2} B(x'x + KI)^{-2} \beta$$

$$= r_{1}(K) + r_{2}(K).$$
 (5.9.b)

puesto que
$$\hat{\beta}^* = z\hat{\beta} = Z(X'X)^{-1}X'Y$$
, se sigue que
$$V(\hat{\beta}^*) = \delta^2 Z(X'X)^{-1}Z'. \tag{5.10}$$

Como r_1 (K) = Traza $\sigma^2 Z(X'X)^{-1} Z$ = suma de las varianzas de todas las $\hat{\beta}_i^*$, diremos que r_1 (K) es la "varianza total" de los es timados de los parámetros.

Significado de r₂(K)

Como $r_2(K) = (Z\beta - \beta)!(Z\beta - \beta)$ y $r_2(0) = 0$, puesto que $Z(0) = I_2(K)$ puede ser considerado como el cuadrado del sesgo al usar $\hat{\beta}^{**}$ en vez de $\hat{\beta}$.

5.3.2 Comportamiento comparado de r₁(K) y r₂(K)

(i) Según (5.9.b) se tiena

$$r_1(K) = \sigma^2 \sum_{i=1}^{k} \frac{\lambda_i}{(\lambda_i + k)^2}$$

lumgo $r_1(K)$ es una función continua, monótona decreciente en K, y :

$$\lim_{k\to 0} r_i(k) = 6^2 \sum_{i=1}^{4} r_i(k) = 0$$
. [5.11.a.]

(ii)
$$\frac{d}{dk} r_i(k) = -26^2 \sum \frac{\lambda_i}{(\lambda_i + k)^3}$$
, [5.11.b]

luego :

$$\lim_{k \to 0^{+} dk} \frac{d}{dk} \gamma(k) = -26^{2} \sum_{k \to 0^{+} dk} \frac{1}{1}$$

$$\lim_{k \to \infty} \frac{dr}{dk} (k) = 0,$$
(5.11.c)

y si X'X deviene en singular, en este caso algún $\lambda_i \to 0$, Y $\lim_{\kappa \to 0^+} \frac{d}{d\kappa} \, (\kappa) \longrightarrow +\infty$

(iii) Según (5.9.6), $r_2(K) = K^2 \beta'(X'X + KI)^{-2}\beta$; afirmamos que: $r_2(K)$ as una función continua monótona creciente.

<u>Prueba</u>. Sea P' la matriz ortogonal que diagonaliza a W², esto es

$$PW^{2}P' = D = \begin{bmatrix} d_{1} & & 0 \\ & \ddots & 0 \\ & & d_{t+1} \end{bmatrix},$$

les dison los valores propios de W, según (5.2.b) se tiene $d_i = \frac{1}{(\lambda_i + K)^2} \cdot \text{Luego},$

$$r_{2}(K) = K^{2}(P\beta)^{i}PW^{2}P^{i}P\beta = K^{2}(P\beta)^{i}Q_{i}$$

$$= K^{2}\sum_{i}(P\beta)^{i}Q_{i}^{2}$$

$$= K^{2}\sum_{i}(X)^{i}Q_{i}^{2}, \qquad (5.12.a)$$

donds \propto = P\$. Cada elemento λ_{i} + K as positivo, por lo tan to r_{p} as continua , para K \neq 0

$$r_2(K) = \sum_{i=1}^{\infty} \frac{\alpha_i^2}{(1+(\frac{\lambda_i}{K}))^2}$$
; $\lambda_i > 0$.

Por lo tanto, r₂ es monótona creciente.

Notemos también que

(iv)
$$\frac{d}{dk} r_2(k) = 2k \sum \frac{\lambda_i \alpha_i}{(\lambda_i + K)^3}$$
 (5.12.c)

luego

$$\lim_{k \to 0^+} \frac{d Y_2}{d k}(k) = 0$$
, y $\lim_{k \to \infty} \frac{d Y_2}{d k}(k) = 0$. (5.12.d)

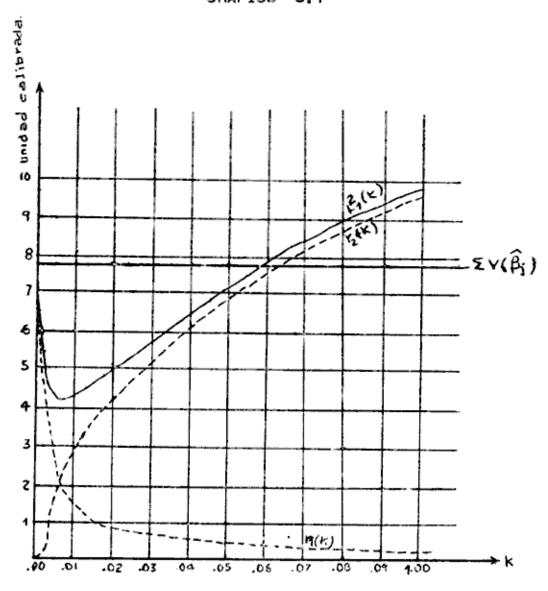
Estos conceptos se visualizan en el gráfico 5.1.

5.3.3 Teoremo de existencia do mínimo de $E(L_1^2(K))$ para K > 0.

Hemos visto que $E(L_1^2(K)) = r_1(K) + r_2(K)$ y el gráfico (5.1) nos muestro un mínimo de $L_1^2(K)$.

Como lim. $r_2(K) = \beta'\beta$, a medida que la magnitud de β crece, se puede pensar que alguna vez $E(L_4^2)$ toma su mínimo en





cero; pero esto no ocurre.

Teorema 5.1 (teorema de existencia). Siempre existe un K > 0 tol que:

$$\varepsilon(L_1^2(K)) < \varepsilon(L_1^2(0)) = 6^2 \sum_{i=1}^{1} (5.13)$$

<u>Prueba</u>. Sogún (5.11,b)_y(5,12,c) se tiene

$$\frac{d}{dk} E(L_1^2(K)) = \frac{d}{dk} \frac{\pi}{\pi} (K) + \frac{d}{dk} \frac{\pi}{\pi} (K)$$

$$= \sum \frac{(K \propto_i^2 - \sigma^2) 2\lambda_i}{(\lambda_i + K)^3}$$

$$\leq (K \propto_{max}^2 - \sigma^2) (\sum \frac{2\lambda_i}{(\lambda_i + K)^3}),$$

 $\frac{d}{dK} \in (L_1^2(K)) \text{ as menor que cero siempre que } K \land \frac{g^2}{g^2_{max}}, \text{ luego}$ $E(L_1^2(K)) \land E(L_1^2(O)).$

5.4. Recomendaciones para la selección de un $\hat{oldsymbol{eta}}^*$ adecuado

Naturalmente el usuario es el que toma la decisión final acerca del sesgo, la varianza y la magnitud del error cua drático medio. Pero, para su ayuda, se den los criterios que deben normer la selección de K.

- (i) Para un determinado valor de K el sistema debe mostrarse estable y los valores propios de X'X + KI deben ser mayores que cero.
- (ii) Los coeficientes de $\hat{m{eta}}_i^*$ deben tener un valor absoluto moderado de acuerdo a la razón de cambio de la variable asociada ${f x}_i$.
- (iii) Los coeficientes que eparentemente tienen signo inco rrecto, deben cambiar de signo al tomar K=0.
- (iv) RSS($\hat{\beta}^{*}$) no debe ser muy grande con respecto a RSS($\hat{\beta}$) ni con respecto a una varianza razonable acorde con el proceso de generación de datos.

RESUMEN

Al analizar un proceso que relaciona variables independientes x_i con una variable dependiente y, los conceptos
técnicos emorgentes del contexto del proceso no son suficientes para explicar la relación que se busce. Es así que se plan
teen modelos, entre ellos los de tipo lineal, cuya discusión
ha sido el objetivo mayor de este trabajo.

Con la hipótesis de que el modelo es lineal, la determinación de la magnitud de la influencia de cada x_i en la
respuesta final y, en la mayoría de los casos se ve afectada
considerablemente por parturbaciones, las cuales pueden surgir
por la inclusión de variables extrañas y/o porque las varia bles de predicción consideradas guardan relación entre si (mul
ticolinealidad) y que tomadas en conjunto dificultan la deter
minación de la influencia de cada una de ellas en la respuesta final. En este sentido las recomendaciones del presente tra
bajo son:

- i) Determinar con claridad el uso al que se destine la regre sión.
- ii) Considerando la variabilidad de los coeficientes β a estimar, la eliminación de variables poco significativas, es ventajosa, no obstarte el esego que se introduce por eliminación.
- ii) En una primera fase, además de los estadísticos R^2 y cociente F, un criterio bastante elaborado que ayuda en la eliminación de variables es $C_{\rm p}$.
- iv) En una segunda fase, el cálculo de los valores propios más cercanos a cero y de sus vectores propios esociados son

de gran eyude en el enálisis y en la eliminación de variables.

v) Finalmente, si la multicolinealidad de los datos persiste y
ya no es recomendable la eliminación de variables, la regra sión por eristas eyude a la determinación de una mejor ecuación
de ajuste .

BIBLIOGRAFIA

- Averson, J.N. y Mc Cobe, G.P. (1973). Yariable selection in a regression analysis. Department of Stat. Univ. of Kentucky.
- Anderson, V.L. y Mc Leen, R.A. (1974). Dessign of experiments

 Mercel Dekker.
- Draper, N.R. y Smith, H(1966). Aplied Regression Analysis
 Wiley, New York.
- Foreythe, G.E.y Moler, C.B.(1967), Computer solution of linear algebraic systems. Prentice Hall.
- Graybill, F. A(1961). An introduction to linear statistical models, Vol.1. Mc Graw Hill.
- Gunst, A.F. y Mason, A.L.(1977). Advantages of examining multicollinearities in Asgression Analysis.

 Biometric, Marzo.
- Hocking, R.R(1976). The analysis and selection of variables in linear regression. Biometrics 32, 1 48
- Hoerl,A.E. y Kennard,R.W.(1970), Ridge regression : Biased estimation for nonortogonal problems.

Technometrics Vol.12, NO 1.

- Huang, D. S(1970). Regression and econometric methods Wiley. New York.
- Johnston, J. (1972) Linear Algebra. Addison Wesley.
- Mallows, C.L (1973). Some comments on $C_{\rm p}$. Technometrics. Vol. 15 Marquardt, P.W.y Snee, R.D (1973). Ridge regression.

Department of stat. Univ. of Kentucky.