

Caja
324

Nº
25

MINISTERIO DE ECONOMIA
SECRETARIA DE COORDINACION Y PLANIFICACION ECONOMICA
INSTITUTO NACIONAL DE PLANIFICACION ECONOMICA
BIBLIOTECA

Hipólito Yrigoyen 250, piso 8º, of. 801/C,
SERIE DE INVESTIGACIONES Nº 12
Buenos Aires (Argentina)

Dros. FERNANDO FERRERO Y CARLOS E. SANCHEZ

MODELO LINEAL PARA LA IDENTIFICACION DE INTERACCIONES

*Presentado: Reunión de Centros de Investigación
en Economía, La Plata, 1969*



UNIVERSIDAD NACIONAL DE CORDOBA
FACULTAD DE CIENCIAS ECONÓMICAS
INSTITUTO DE ECONOMIA Y FINANZAS
1971

MODELO LINEAL PARA LA IDENTIFICACION DE INTERACCIONES *

INTRODUCCION

El presente trabajo tiene por objeto desarrollar un método para identificar y estimar interacciones en un modelo lineal donde todas las variables son de naturaleza dicotómica. Lejos de ser ésta una simple cuestión de transformación de variables, en el sentido de que variables originariamente cuantitativas se transforman en cualitativas, el problema central que aquí se plantea es lograr un conjunto de variables en el que cada una de ellas controla todos los niveles posibles de las restantes, facilitando de esta manera la estimación del efecto puro o aislado que a cada una de ellas le corresponderá. Tales efectos, por el método de construcción adoptado, no dependerán en modo alguno de los valores o niveles que asuman las restantes variables.

Las variables que en este modelo se utilizan no son otra cosa que las que resultan de la conjunción de los diversos niveles de todas las variables y en la medida que la estimación puntual y la docimasia de hipótesis revelen la presencia de coeficientes significativos, tales combinaciones estarán denunciando un comportamiento diferencial por parte de las observaciones o individuos que las componen. Cuando una categoría cualquiera tiene precisamente estas características se la denomina una interacción.

* Serie de Investigaciones del Instituto de Economía y Finanzas N° 12.

Conviene aclarar que el propósito del trabajo es esencialmente metodológico y aun cuando se ha tratado de probarlo empíricamente para investigar la posesión de un bien durable, en este caso automóviles, no se pretende con ello elaborar un marco teórico o conceptual de tal fenómeno.

ESTRUCTURA DEL MODELO

En el modelo lineal que analizaremos, que no es de regresión ni de varianza ni de covarianza, todas las variables intervinientes tanto exógenas como endógenas son de naturaleza dicotómica. Una determinada variable asume el valor uno o cero según que el individuo u observación pertenezca a una cierta categoría, o posea o no un determinado atributo. En este caso particular relacionamos la posesión de un cierto bien durable —automóviles— con un conjunto de variables, tales como ingreso familiar, edad del jefe, posición ocupacional, educación, etc., con lo que el modelo por ser lineal adoptará la forma

$$Y = XB + u, \text{ (sin formular por el momento supuesto alguno acerca de la variable aleatoria } u \text{)}$$

y en términos de matrices y vectores, la siguiente composición:

Posee auto-móvil	Ingresos			Edad		Ocupación del jefe			...
	a_1-a_1	a_2-a_2	a_4-a_3	b_1-b_1	b_3-b_2	Asalariado	Inter-medio	Jefe alto nivel	
0	1	0	0	1	0	0	1	0	...
1	0	1	0	0	1	0	0	1	...
.
.

De este modo el individuo ubicado en la primera fila se caracteriza por pertenecer al primer tramo de ingresos, al primer inter-

valo de edad y al segundo de ocupación y no posee automóvil. El segundo sí posee automóvil, pertenece al segundo tramo de ingresos, al segundo de edad y al tercero de ocupación. Se tendrán de esta forma tantas filas cuantas observaciones existan y un número de columnas coincidente con el número de variables que se analiza.

En forma abreviada,

$$Y_i = \begin{cases} 1 & \text{si el } i\text{-ésimo individuo posee} \\ & \text{automóvil} \\ 0 & \text{en caso contrario} \end{cases}$$

$$X_{ij} = \begin{cases} 1 & \text{si el } i\text{-ésimo individuo pertenece a la} \\ & \text{j-ésima categoría} \\ 0 & \text{en caso contrario} \end{cases}$$

Dado un conjunto de observaciones del tipo señalado, ha de existir ciertamente alguna relación entre las variables independientes y la variable que convencionalmente designamos como dependiente. Es necesario encontrar, en otras palabras, algo que nos explique por qué una mayor proporción de los individuos que pertenecen a X_i poseen automóvil, en tanto que los que pertenecen a X_j muestran una frecuencia menor en orden a la posesión de automóviles.

Puesto que las variables que en realidad se utilizarán son distintas a las indicadas en la página anterior y para evitar la confusión entre variables y sub-variables, se usará la notación: $X_{ij}^{k,m}$ para representar la siguiente categoría: el primer supra índice indicará una de las variables originales; ingreso, educación, edad, etc.; el segundo supra índice indicará el número de niveles en que tal variable ha sido dividida o también el número de variables dicotómicas que da origen; el primer sub-índice individualizará al i -ésimo individuo; y el segundo, el nivel específico al que el individuo pertenece. De esta forma, en el ejemplo anterior la primera columna se representará por $X_{i1}^{1,2}$; la segunda, $X_{i2}^{1,2}$; la cuarta, $X_{i4}^{2,2}$; la séptima, $X_{i7}^{3,2}$, etc. A su vez, un elemento cualquiera de la matriz.

$$X_{ij}^{k,m_k} = \begin{cases} 1 & \text{si el } i\text{-ésimo individuo pertenece al } j\text{-ésimo nivel de la } k\text{-ésima variable (ésta ha sido dividida en } m_k \text{ niveles)} \\ 0 & \text{en caso contrario} \end{cases}$$

De hecho que si $X_{ii}^{k,m_k} = 1$ ello implicará que $X_{ij}^{k,m_k} = 0$ para todo $j \neq i$. Esta particular característica hace que la matriz sea de rango incompleto.

Es necesario destacar que si bien las variables que aquí aparecen son de naturaleza dicotómica —lo que restringiría el uso del método sólo a variables no cuantitativas—, éstas en su forma original pueden ser de cualquier tipo. Si el nivel de medición original es típicamente el de una variable cuantitativa, como en el caso de Ingresos del ejemplo anterior, siempre será posible dividirla en intervalos y asignar a cada uno de ellos una variable dicotómica específica. De aquí que en lo sucesivo, al referirnos a variables, estaremos más bien implicando categorías o niveles de una cierta variable.

ANÁLISIS DE RELACIONES

A primera vista podría aparecer acertado plantear un modelo de regresión del tipo $Y = XB + u$ y estimar el vector B con las correspondientes dójimas de hipótesis asociadas a B o a un subconjunto de sus elementos. Los coeficientes así estimados darían una medida del efecto aislado de cada variable cuando las restantes permanecen constantes. Tal procedimiento, empero, no es aconsejable ya que, 1) la matriz X es de rango inferior a $\sum m_k$ (número de variables incluidas), lo cual obligará a prescindir de algunas variables o estimar combinaciones lineales de parámetros; 2) en la propia construcción del modelo de regresión lineal se parte del supuesto de que las variables sólo tienen efectos aditivos y de allí que el coeficiente de regresión estimado sea una medida apropiada del efecto aislado de cada variable. Sin embargo, el método de regresión múltiple es primordialmente eficaz en orden a la predicción de la variable dependiente, pero sí lo que se busca es la deter-

minación de los parámetros de una relación estructural, la regresión múltiple no es en tal circunstancia el método más apropiado¹; 3) cuando se combinan un grupo de variables económicas y demográficas el supuesto de aditividad se hace más insostenible; no es fácil aceptar que el ingreso de una persona pueda ser disociado de su posición ocupacional o de su edad y más aún admitir que cuando estas variables actúan aisladamente sus efectos son semejantes a cuando lo hacen conjuntamente. Cuando se maneja un conjunto de variables que presentan marcadas interacciones, carece de sentido investigar el efecto de cada una de ellas consideradas aisladamente².

Todas estas circunstancias aconsejan utilizar un modelo donde las variables intervinientes resultan de una conjunción de los diversos niveles de todas las variables. De esta forma las nuevas variables así construidas representarán el efecto conjunto y simultáneo de cada nivel de las restantes variables y podrán en consecuencia indicar un efecto puro o aislado que en modo alguno dependerá de si las otras variables permanecen o no constantes.

A las variables primitivas, que por lo expuesto conllevan los efectos confundidos de las restantes, las denominaremos "genéricas" o "complejas"; a las nuevas variables, que representan efectos exclusivos y atribuibles a ellas mismas, las denominaremos "simples" o "puras".

Para pasar de las complejas a las puras definiremos la operación de producto lógico de categorías que se simbolizará por (x), de forma que las nuevas variables serán

$$Z = X^{1,m_1} (x) X^{2,m_2} \dots (x) X^{K,m_K} \quad \text{(siendo K el número de variables complejas)}$$

$$Z = X^{1,m_1} (x) X^{2,m_2} \dots (x) X^{K,m_K}$$

1. TITNER, G.: *Econometrics*, New York, John Wiley and Sons; 1952.

2. MORGAN, J.H. and SONQUIST, J.A.: "Some results from a nonsymmetrical branching process that looks for interaccion effects". *Proceedings of the Social Statistics Section of the American Statistics Association*, 1963, 40-45.

donde Z_{j_1} representará la combinación de los primeros niveles de todas las variables complejas; Z_{j_2} la combinación de los primeros niveles de $K-1$ y el segundo de la K -ésima, etc., en general,

$$Z_{j_1} = X_{j_1}^{1,m_1} (x) X_{j_2}^{2,m_2} \dots (x) X_{j_k}^{K,m_k}$$

siendo $j_1 = 1, 2, \dots, m_1$

$j_2 = 1, 2, \dots, m_2$

.....

$j_k = 1, 2, \dots, m_k$

Con estas nuevas variables la representación del modelo cambia en el sentido de que ahora un individuo sólo puede pertenecer a una única variable Z_{j_1} ; la intersección entre Z_{j_1} y Z_{j_2} , siendo $j_1 \neq j_2$, es nula ya que una de ellas incluye algún nivel que la otra excluye. Obviamente a una cierta Z_{j_1} podrán pertenecer más de un individuo.

La matriz del modelo tendrá en consecuencia la siguiente composición:

Y	Z_{j_1}	Z_{j_2}	Z_{j_3}	Z_{j_m}		Z_{j_n}
1	0	1	0	...	0	0
0	1	0	0	...	0	0
0	1	0	0	...	0	0
.
1	0	0	1	...	0	0
r	n_1	n_2	n_3	...	n_j	n_m

en la cual r es el número de individuos que en la muestra poseen automóvil.

n_j es el número de individuos que pertenecen a Z_{j_1}

$$y \quad n = \sum_{j=1}^m n_j \quad y \quad m = \sum_{k=1}^K m_k$$

ESTIMACION DE PARAMETROS

Para la variable j , el producto vectorial $Y'Z_j$ es el número de individuos que del total n_j de la categoría poseen automóvil, sea $p_j = Y'Z_j/n_j$ la proporción que poseen automóvil en la categoría j . A su vez la matriz Z tiene la característica de que $Z'Z = I \cdot n$, (I es una matriz unidad de $m \times m$ y $n' = (n_1, n_2, \dots, n_m)$) y por consiguiente su inversa será

$$(Z'Z)^{-1} = I \cdot n^{-1},$$

El producto $Z'Y = n \cdot p$, donde $(n \cdot p)' = (n_1 p_1, n_2 p_2, \dots, n_m p_m)$.

Los estimadores mínimo cuadráticos del modelo serán por tanto $B = (Z'Z)^{-1}Z'Y = I \cdot n^{-1} n \cdot p = p$, ($p' = p_1, p_2, \dots, p_m$).

De forma tal que por la propia estructura del modelo que estamos considerando los estimadores mínimo cuadráticos son directamente iguales a las frecuencias marginales que indican la proporción de automóviles que se poseen en cada categoría. Tal circunstancia facilita enormemente el cálculo y permite realizarlos cualquiera sea el tamaño de la matriz.

Al igual que se definió la operación de intersección, se puede definir la operación de unión o suma lógica de dos categorías, implicando con ello una nueva categoría, menos pura que sus componentes, y a la cual pertenecerán los individuos que correspondían a una u otra categoría. Así por ejemplo $Z_{j+j'} = Z_j (+) Z_{j'}$ constará de $n_j + n_{j'}$ individuos y de ellos $Y'(Z_j + Z_{j'})$ poseerán automóviles. Se puede comprobar que por la naturaleza de la matriz Z , el coeficiente estimado de esta nueva categoría es

$$p_{j+j'} = (n_j p_j + n_{j'} p_{j'}) / (n_j + n_{j'})$$

que se puede extender a cualquier número de categorías. Además en base a esta operación se podrá más adelante estimar la contribución a la variación explicada de un subconjunto cualquiera de variables.

ANÁLISIS DE VARIANZA

Los coeficientes estimados oscilan entre cero y uno; si todos los individuos que pertenecen a una cierta categoría poseen automóvil, el valor de p será uno; si ninguno posee automóvil, el valor de p será 0. En base a las relaciones definidas, se puede comprobar que la variación total, $y'y = r - r^2/n$, la suma residual de cuadrados, es

$$e'e = Y'Y - B'Z'Y = r - p'(n.p.) = r - \sum_{i=1}^m n_i p_i^2,$$

la variación explicada,

$$= \sum_{i=1}^m n_i p_i^2 - r^2/n$$

el coeficiente de correlación múltiple,

$$R^2 = (\sum n_i p_i^2 - r^2/n) / (r - r^2/n)$$

la varianza de los coeficientes estimados,

$$\text{var}(B_i) = (r - \sum_{i=1}^m n_i p_i^2) / n_i (n - m)$$

$$\text{cov}(B_j, B_p) = 0 \quad \text{para todo } j \neq p$$

Por su parte, la variación explicada puede descomponerse del siguiente modo:

$$\begin{aligned} \sum n_i p_i^2 - r^2/n &= n_1 p_1^2 - (n_1 p_1)^2/n + \sum_{i=2}^m n_i p_i^2 - (\sum_{i=2}^m n_i p_i)^2/n - \\ &\quad - 2(n_1 p_1 \sum_{i=2}^m n_i p_i) / n \\ &= \sum_{i=2}^m n_i p_i^2 - r_1^2 / (n - n_1) + n n_1 (p_1 - r/n)^2 / n - n_1 \end{aligned}$$

en la que el primer sumando representa la variación explicada cuando no se incluye en el modelo la categoría Z_1 (siendo $r_1 = \sum_{i=2}^m n_i p_i$), y el segundo la adición debida a la incorporación de tal categoría.

MODELO LINEAL PARA LA IDENTIFICACION DE INTERACCIONES

Formalmente la tabla de análisis de varianza adoptará la siguiente disposición:

Fuente de variación	Suma de cuadrados	Grados de lib.
Debido a Z_1	$nn_1 (p_1 - r/n)^2 / (n - n_1)$	1
Debido a Z_1, Z_2, \dots, Z_m , excepto Z_1	$\sum_{j=2}^m n_j p_j^2 - r^2 / (n - n_1)$	$m - 2$
Variación explicada	$\sum_{j=1}^m n_j p_j^2 - r^2/n$	$m - 1$
Variación residual	$r - \sum_{j=1}^m n_j p_j^2$	$n - m$
Variación total	$r - r^2/n$	$n - 1$

Vale la pena destacar que en el problema usual de la regresión lineal un coeficiente que no difiere significativamente de cero conduce por lo general a la afirmación de que tal variable no explica o no contribuye a la variación explicada por la regresión. Dejando de lado los reparos que tal generalización merece, la regla comúnmente aceptada es de que tal variable puede ser excluida del modelo sin que por ello éste pierda su valor predictivo. En el caso presente, sin embargo, un coeficiente que se aproxima a cero o a uno (supuesto que r dista de ambos extremos) es en cualquier caso significativo, bien sea porque está revelando una asociación positiva con la variable dependiente o bien porque está indicando una correlación negativa.

El valor de r representa en esencia un promedio y una categoría cualquiera es significativa en la medida que su coeficiente estimado se aparta del valor promedio; circunstancia ésta que evidenciará un comportamiento diferencial por parte de los miembros que pertenecen a esa categoría y existirá en consecuencia una interacción en la combinación de niveles que tal categoría supone. La detección de interacciones es precisamente el objetivo del método que desarrollamos.

Si guiendo un razonamiento similar, si lo que interesa es la variación atribuible a dos categorías, que representan desde luego un nivel de interacción más general, tal contribución se podrá estimar a través del valor de

$$(n_1 + n_2) n \left(\frac{r}{n} - \frac{(n_1 p_1 + n_2 p_2)}{(n_1 + n_2)} \right) / (n - n_1 - n_2)$$

DOCIMASIA DE HIPOTESIS

Para no sujetar la validez de las conclusiones a los supuestos que generalmente acompañan los modelos lineales —supuestos cuya inclusión es necesaria a los efectos de docimar hipótesis referentes a los parámetros— se ha optado por un test no paramétrico que se basará en los coeficientes estimados. Bajo los supuestos de la hipótesis nula que afirmarán la no significatividad de las categorías incluidas, el test o estrategia radicarán en la distribución del escalar $Y'Z_1$, número de personas que del total n_1 de la categoría Z_1 poseen automóvil. Cualquiera sea la forma o los parámetros de la distribución de donde provienen las observaciones, el estadístico utilizado tendrá una distribución hipergeométrica de parámetro r/n .

El vector Y es una muestra de tamaño n en la cual cada individuo ha sido clasificado dicotómicamente conforme a la posesión de automóvil; el vector Z_1 es una submuestra de la anterior, de tamaño n_1 , y en la que cada individuo ha sido a su vez reclasificado según su pertenencia a dicha categoría. En cualquier caso $n_1 \neq n$.

Podemos considerar a Z_1 como una submuestra de Y o bien condicionar Z_1 por lo observado en Y , en este último caso Z_1 equivaldría a una muestra de la población que Y representa. Ciertamente lo ideal sería conocer R/N (siendo $E(r/n) = R/N$) que no es otra cosa que el parámetro de la población, pero la cuestión esencialmente no varía si se enfoca el problema combinatorio correctamente.

Consideremos primeramente un caso abreviado en el que intervienen 20 observaciones y tres variables y que presentan la siguiente disposición:

MODELO LINEAL PARA LA IDENTIFICACION DE INTERACCIONES

Obs. n°	Y	Z ₁	Z ₂	Z ₃
1	1	0	1	0
2	0	0	0	1
3	0	0	0	1
4	1	0	1	0
5	0	0	0	1
6	0	1	0	0
7	1	1	0	0
8	1	0	1	0
9	0	0	0	1
10	1	0	1	0
11	0	0	0	1
12	1	0	1	0
13	0	1	0	0
14	1	0	1	0
15	1	1	0	0
16	0	0	0	1
17	0	0	0	1
18	0	0	0	1
19	0	0	0	1
20	0	0	0	1

$$r = 8 \quad Z'_1 Y = 2 \quad Z'_2 Y = 6 \quad Z'_3 Y = 0$$

$$n = 20 \quad n_1 = 4 \quad n_2 = 6 \quad n_3 = 10$$

$$p_1 = 0.5 \quad p_2 = 1 \quad p_3 = 0$$

Analizando los resultados de la variable que representa Z₃ cuyo coeficiente es p₃ = 0, podría parecer a primera vista que a pesar de ser tan bajo, tal resultado es compatible con una muestra de tamaño 10 proveniente de una población en la que 40% de los individuos tienen automóvil, o equivalentemente de una población de 20 individuos, 8 de los cuales poseen automóvil. Un resultado tan bajo como el observado podría ser consecuencia de fluctuaciones muestrales, por lo que la hipótesis nula establecerá que tal muestra proviene en efecto de una población compuesta por 20 individuos y en la que un 40% de ellos poseen automóvil. A los fines de docimar tal hipó-

tesis utilizaremos el estadístico $Z_1 Y$ cuya distribución bajo la hipótesis nula responde a una hipergeométrica de parámetro 0,40; la distribución de $Z_1 Y$ será por consiguiente:

$$P(Z_1 Y=0/H_0) = \Pr (p_1=0) = \binom{8}{0} \binom{12}{10} / \binom{20}{10} = 0.00036$$

$$P(Z_1 Y=1/H_0) = P (p_1=0.1) = \binom{8}{1} \binom{12}{9} / \binom{20}{10} = 0.00953$$

$$P(Z_1 Y=2/H_0) = P (p_1=0.2) = \binom{8}{2} \binom{12}{8} / \binom{20}{10} = 0.07502$$

.....

Si los individuos hubieran sido efectivamente extraídos de una población en la que 8 de cada 20 poseen automóvil, la evidencia muestral nos colocaría frente a la siguiente alternativa, a) si aceptamos la hipótesis nula, deberemos por fuerza admitir que un evento de probabilidad tan baja como 0.00036 ha podido ocurrir, b) o bien rechazamos la hipótesis nula fundándonos en que el valor observado queda comprendido dentro de la zona de rechazo. A un nivel de significación del 5% la misma decisión cabrá adoptar si en vez de observar el valor $p_1=0$ hubiéramos observado $p_1=0,1$.

Por lo que respecta a la variable Z_1 si se calculan las probabilidades para cada uno de los valores que puede asumir el estadístico $Z_1 Y$ se encontrará que el valor observado $p_1=0.50$ no motivará el rechazo de H_0 . Para la variable Z_2 , un valor tan alto como el observado sólo tiene una probabilidad de 0,0072 de provenir de una población con parámetro 0,40. De un modo general, el test consistirá en establecer la siguiente hipótesis:

$$H_0 : p_1 = r/n$$

$$H_1 : p_1 \neq r/n$$

basándose en la distribución de $n_1 p_1$ que para muestras pequeñas tiene una distribución exacta,

MODELO LINEAL PARA LA IDENTIFICACION DE INTERACCIONES

$$n_j p_j \sim \begin{pmatrix} r \\ n_j p_j \end{pmatrix} \begin{pmatrix} n-r \\ n_j - n_j p_j \end{pmatrix} / \begin{pmatrix} n \\ n_j \end{pmatrix}$$

y para grandes muestras su aproximación normal,

$$n_j p_j \sim (\sqrt{2\pi h})^{-1} \exp(-1/2 \cdot (n_j p_j - n_j r/n)^2 / h^2)$$

siendo $h = (n_j(r/n) (1-r/n))^{1/2}$

APLICACION A UN CASO CONCRETO

En base a las encuestas que sobre empleo y desempleo realiza periódicamente el Instituto de Economía y Finanzas se han obtenido datos sobre posesión de bienes durables. A los fines del presente trabajo se adoptó la correspondiente a abril de 1967 y allí se extrajo información sobre 1364 unidades familiares. Las variables que se utilizaron fueron:

A) Ingresos del grupo familiar

- A.1 menos de \$ 30.000 mens.
- A.2 30.000 - 70.000.—
- A.3 más de 70.000.—

B) Edad del Jefe

- B.1 30-39 años
- B.2 40-49 años
- B.3 50-59 años
- B.4 60 y más

C) Categoría ocupacional del Jefe

- C.1 No calificados (tareas manuales y no manuales)
- C.2 Calificados (tareas manuales y no manuales)
- C.3 Técnicos, profesionales, jefes y directivos

E) Educación del Jefe

- E.1 Sin ningún tipo de enseñanza
- E.2 Enseñanza primaria (completa o incompleta)

E.3 Ciclos medios (completos e incompletos) y universitario incompleto.

E.4 Universitario completo

Y) Variable dependiente

0 No posee automóvil

1 Posee automóvil

La tabla I presenta los resultados que se obtuvieron por aplicación del método descrito. Para una explicación mejor de la variable dependiente hubiera sido necesario incorporar otras variables, de forma tal que los coeficientes estimados de cada celda se aproximaran a cero o uno. Considérese, por ejemplo, la variable A_3, B_3, E_4, C_3 en la que aparecen 11 individuos, de los cuales 8 poseen automóvil y 3 no poseen. Es fácil conjeturar que las ocho unidades familiares que poseen automóvil tienen que ser diferentes a las restantes y que tales diferencias pueden revelarse a través de alguna variable no incluida en el modelo. O tal vez las mismas variables utilizadas podrían indicar tales diferencias si es que no se hubiera recurrido a un nivel de agregación (o de formación de niveles) tan alto como el que aquí se usó. El criterio que sirvió de base para la formación de niveles fue un tanto arbitrario; alternativamente y desde luego con mejores resultados podría utilizarse un método de construcción de niveles que maximice la homogeneidad interna de cada celda (o de cada variable), entendida ésta en el sentido que cada celda contenga el mayor número posible de individuos con automóvil y por consiguiente el menor número posible de individuos sin tal atributo. Como caso límite las celdas deberían formarse de modo que algunas incluyan únicamente unidades que poseen automóviles y las restantes celdas con todos los individuos que no poseen automóvil.

Las ventajas de este procedimiento pueden apreciarse en las fórmulas ya desarrolladas. La suma residual de cuadrados

$$e'e = r - \sum_{j=1}^k np_j^2$$

MODELO LINEAL PARA LA IDENTIFICACION DE INTERACCIONES

tiene la propiedad de que si todos los coeficientes son iguales a cero o uno,

$$\sum n_j p_j^2 = \sum n_j p_j = r$$

y en consecuencia

$$e'e = 0.$$

El análisis de varianza de la tabla I dio los siguientes resultados:

Fuente de variación	Suma de cuadrados	G.L.
Variación explicada	73.070	144
Variación residual	169.185	1219
Variación total	242.255	1363

Dentro del esquema desarrollado se puede estimar la contribución a la variación explicada proveniente de las variables que en la tabla I aparecen con coeficientes significativamente altos. Para ello, sólo basta unir todas las categorías excepto aquellas cuya contribución se quiere determinar, obteniendo una nueva variable cuyo coeficiente será

$$p'_j = (\sum n_j p_j) / \sum n_j$$

(donde el subíndice j se extiende a todas las categorías excepto las significativas)

$$= 169/965 = 0,1751$$

Con este nuevo modelo, la tabla de análisis de varianza se convierte

Fuente de variación	Suma de cuadrados	G. L.
Variación explicada	69.950	44
Variación residual	172.305	1319
Variación total	242.255	1363

Para estimar la contribución de p'_j

$$n \cdot n'_j (p'_j - r/n)^2 / n - n'_j = 1364.965 (0.2309 - 0.1751)^2 / 399 = 10,263$$

Luego la contribución de las restantes variables es de aproximadamente 85% de la variación explicada.

Para docimar individualmente cada coeficiente se construyó una tabla (Tabla II del Apéndice) con los valores críticos del 90 y 95% para muestras de tamaño inferior o igual a 30. Cuando los valores de n_i superaban tal límite, se utilizó la aproximación normal ya descripta.

INTERPRETACION DE LOS RESULTADOS

Del análisis de los coeficientes estimados (p_i) y de sus niveles de significación se desprende que, en primer lugar, todas las variables presentan una correlación positiva entre posesión de automóvil y su nivel, o sea que a medida que se asciende en el nivel de ésta, mayor es el valor del coeficiente. En el caso de edad, sin embargo, esto es observable sólo en dos tramos, menos de 40 años en el que la posesión evidencia un comportamiento significativamente bajo y, superada esa edad crítica, un tramo en el que los coeficientes son aproximadamente iguales para cualquier nivel de edad.

Con respecto a la contribución de cada categoría a la explicación total del modelo (30%), se observa que un 85% de esa variación explicada está dada por los extremos de las variables, es decir, alta educación, alto nivel de ingresos, alta índole ocupacional, etc. y también por ciertos niveles intermedios de estas variables que se interaccionan de un modo particular.

La no posesión de automóviles, por su parte, es explicada en su totalidad por el nivel más bajo de ingreso y los dos primeros niveles de educación, cualquiera sea la índole ocupacional y edad.

Finalmente, existe una zona en la que los coeficientes estimados no difieren significativamente del promedio de todas las categorías (0,2309). Esta zona de comportamiento no diferencial está primordialmente encerrada entre el tercer nivel de educación con ingresos bajos y el segundo nivel de educación con ingresos intermedios.

Las variables incluidas permiten formular algunas consideraciones sobre el tema, siguiendo los lineamientos teóricos dados en un

trabajo anterior (3) en el que se analiza la conducta del consumidor con respecto a la posesión de bienes de consumo durables.

Mediante el análisis de regresión se vincula allí el stock de bienes durables poseído por el consumidor con las variables ingreso corriente y edad. No se utiliza, sin embargo, una regresión múltiple que sólo permitiría medir el efecto de cada variable sobre el stock de durables poseído, aislando la influencia de la otra pero sin controlar el nivel al cual se la mantiene constante. Por el contrario, se intenta controlar a cada variable para niveles dados de la otra en base a regresiones simples en las que, para cada uno de los distintos niveles definidos, tanto para ingreso como para edad, se relaciona el stock con la variable "libre".

Partiendo del supuesto de que el transcurso del tiempo es un requisito necesario para la acumulación de riqueza, se deduce que las variaciones de edad, conjugadas con distintos niveles de ingreso corriente, dan lugar a otros tantos niveles de riqueza.

Esta variable, a su vez, es definida como el patrimonio neto del consumidor en el que se incluye la capacitación del individuo o riqueza humana. Este último aspecto sin embargo no es controlado dentro del análisis, ya que ni la variable ingreso, ni la variable edad, son "depuradas" del concepto capacitación.

Al incluir ahora las variables educación e índole ocupacional, la medición de la riqueza puede por lo tanto ser refinada un poco más a efecto de comprobar la hipótesis de que la riqueza acumulada, o capitalización del consumidor, actúa como factor determinante de la propiedad del bien durable automóvil.

En base al cuadro posterior pueden identificarse cuatro zonas diferentes. La primera, tramo 1, corresponde a una situación en la que, ingresos y educación bajos, impiden una acumulación de riqueza tal que permita tener acceso al bien automóvil.

La segunda, tramos 2 y 3, incluye a aquellas unidades de consumo que presentan una conducta diferencial con respecto al promedio. A diferencia del caso anterior, un mayor ingreso o un superior

3. SÁNCHEZ, C.E.: "Posesión de Bienes de Consumo Durables", *Revista de Economía y Estadística*, N° 1-2, primero y segundo trimestres de 1968; 37-39.

INSTITUTO DE ECONOMIA Y FINANZAS

nivel de educación, conjugados con los distintos niveles de edad y de índole ocupacional, definen niveles de riqueza en los que las decisiones sobre posesión resultan positivas.

En la tercera zona, tramos 4 y 5, los niveles superiores o de educación hacen que las variaciones de edad o de ocupación produzcan una conducta diferencial, es decir, la mayor riqueza que esto implica lleva a frecuencias significativas y altas.

ANÁLISIS DE LOS RESULTADOS

Tramo	Variables Controladas	Nivel	Comportamiento de las Variables Edad y Ocupación
1	Ingreso	bajo	para cualquier nivel de edad u ocupación
	Educación	bajo	los coeficientes son significativos y bajos.
2	Ingreso	bajo	para cualquier nivel de edad u ocupación
	Educación	medio y alto	los coeficientes no son significativos (igual al promedio).
3	Ingreso	medio	ídem anterior
	Educación	bajo	
4	Ingreso	alto	coeficientes no significativos excepto
	Educación	bajo	para edad y ocupación altas que presentan coeficientes significativos y altos.
5	Ingreso	medio	los coeficientes de edad y ocupación al-
	Educación	medio y alto	tas pasan a ser significativos y altos
6	Ingreso	alto	para cualquier nivel de edad u ocupación
	Educación	alto	los coeficientes son significativos y altos.

Por último, en el tramo 6, con ingreso y con educación altos, las variables edad y ocupación vuelven a perder poder explicativo, esta

vez en forma opuesta al tramo 1. Es decir el consumidor ha podido ajustar su conducta más rápidamente que en los tramos 4 y 5 y en consecuencia, tanto las variaciones de edad como las de la índole ocupacional no producen un comportamiento diferencial.

Estas conclusiones corroboran la hipótesis formulada y sugieren que el método propuesto constituye un eficaz instrumento de análisis. Como ya se señalara anteriormente, la ampliación del marco conceptual permitirá hacer explícita la interacción de variables aquí no consideradas, aumentando el poder explicativo del modelo y brindando en consecuencia una mejor base empírica para la formulación de una teoría sistemática.

TABLA I

	E ₁				E ₂				total
	C ₁	C ₂	C ₃	total	C ₁	C ₂	C ₃	total	
A ₁	B ₁	0 0.000 (-)	0 0.000 (-)	0 0.000 (-)	0 0.000 (-)	1 0.025 (-)	0 0.029 (-)	0 0.027 (-)	2 72
	B ₂	5 0.000 (-)	2 0.000 (-)	0 0.000 (-)	7 0.000 (-)	39 0.087 (-)	33 0.136 (-)	0 0.121 (-)	72 124
A ₂	B ₃	0 0.000 (-)	1 0.250 (-)	0 0.083 (o)	1 0.083 (o)	7 0.097 (-)	10 0.133 (-)	1 0.250 (o)	18 123
	B ₄	9 0.000 (-)	3 0.000 (-)	0 0.023 (-)	12 0.023 (-)	65 0.111 (-)	65 0.094 (-)	3 1.000 (o)	123 235
total	B ₁	1 0.033 (-)	0 0.000 (-)	0 0.000 (-)	1 0.023 (-)	21 0.111 (-)	7 0.094 (-)	1 1.000 (o)	29 235
	B ₂	33 0.018 (-)	10 0.059 (-)	0 0.028 (-)	43 0.028 (-)	168 0.095 (-)	67 0.110 (-)	0 0.375 (o)	235 564

Cada celda constituye una variable. La cantidad en el margen superior izquierdo es el número de personas que poseen automóvil; la del margen inferior izquierdo, el número de individuos que no poseen automóvil. A la derecha de ambas, el valor del coeficiente estimado. Los símbolos (-), (o), (+), indican que el coeficiente es significativamente bajo, no significativo y significativamente alto, respectivamente, al nivel del 95%.

TABLA I (Continuación)

	E ₁				E ₂				total
	C ₁	C ₂	C ₃	total	C ₁	C ₂	C ₃	total	
B ₁	0	0	0	0	0	0	0	0	0
B ₂	0	0	0	0	0	0.000 (e)	0	0.000 (e)	0.000 (e)
B ₃	0	0	0	0	0	0.000	1	0.000 (e)	2 (e)
B ₄	0	0	0	0	0	0.000	1	0.500 (e)	1 (e)
t o t a l	0	0	0	0	0	0.667 (e)	2	1.000 (+)	5 0.833 (+)
t o t a l	0	0	0	0	0	0.000 (-)	3	0.750 (+)	6 0.500 (+)
t o t a l	1	2	0	3	45	0.110 (-)	63	0.178 (-)	123 0.184 (-)
t o t a l	54	19	0	73	362	291	16	669	

TABLA I (Continuación)

	E ₃						E ₄						
	C ₁	C ₂	C ₃	total	C ₁	C ₂	C ₃	total	C ₁	C ₂	C ₃	total	
B ₁	2	0.500 (o)	0	3	0	0	0	0.158 (o)	0	0	2	2	0.182 (o)
	2	14	0	16	0	0	0	16	0	0	2	2	27
B ₂	5	12	7	24	0	1	6	0.436 (+)	1	1.000 (o)	0.667 (+)	7	0.700 (+)
	6	23	2	31	0	0	3	0.778 (+)	0	0	3	3	0.421 (+)
B ₃	4	11	6	21	0	1	4	0.467 (+)	1	1.000 (o)	0.571 (+)	5	0.560 (+)
	3	19	2	24	0	0	3	0.750 (+)	0	0	3	3	56 (+)
B ₄	2	12	5	19	1	0	7	0.513 (+)	1	0.000 (o)	0.636 (+)	8	0.615 (+)
	6	11	1	18	0	1	4	0.833 (+)	0	0	5	5	0.438 (+)
t o t a l	13	36	18	67	1	2	19	0.429 (+)	1	0.667 (o)	0.613 (+)	22	0.628 (+)
	17	67	5	89	0	1	12	0.783 (+)	0	0	13	13	0.410 (+)

TABLA I (Continuación)

	B ₃				B ₄				total
	C ₁	C ₂	C ₃	total	C ₁	C ₂	C ₃	total	
B ₁	1 0.500 (o)	2 1.000 (+)	0	3 0.750 (+)	0	0	0	0	3 0.750 (+)
B ₂	0	3 1.000 (+)	3 0.750 (+)	6 0.857 (+)	0	1 1.000	11 1.000 (+)	12 12.000 (+)	18 0.857 (+)
B ₃	2 1.000 (+)	5 0.714 (+)	5 1.000 (+)	12 0.857 (+)	0	1 1.000	8 0.727 (+)	9 0.750 (+)	22 0.733 (+)
B ₄	0	2	0	2	0	0	3	3	8
A ₃	2 0.667 (o)	2 0.667 (o)	4 1.000 (+)	10 0.833 (+)	0	1 1.000	7 1.000 (+)	8 1.000 (+)	23 0.885 (+)
t o t a l	5 0.714 (+)	12 0.800 (+)	12 0.923 (+)	29 0.837 (+)	0	3 1.000 (+)	26 0.897 (+)	29 0.906 (+)	64 0.810 (+)
B ₁	2	3	1	6	0	0	3	3	15
t o t a l	34 0.233 (o)	65 0.293 (+)	37 0.711 (+)	138 0.327 (+)	2 0.400 (o)	5 0.833 (+)	46 0.708 (+)	53 0.697 (+)	315 0.2309 (+)
B ₁	112	157	15	284	3	1	19	23	1049

TABLA II

VALORES CRITICOS DEL TEST HIPERGEOMETRICO

n_j	95%		90%	
	valor crítico inferior	valor crítico superior	valor crítico inferior	valor crítico superior
2				2
3		3		3
4		3		3
5		4		3
6		4		4
7		5		4
8		5		4
9		5	0	5
10		6	0	5
11		6	0	5
12	0	6	0	6
13	0	7	0	6
14	0	7	0	6
15	0	7	1	7
16	0	8	1	7
17	0	8	1	7
18	0	8	1	7
19	1	8	1	8
20	1	9	1	8
21	1	9	2	8
22	1	9	2	9
23	1	10	2	9
24	1	10	2	9
25	2	10	2	9
26	2	10	2	10
27	2	10	3	11
28	2	11	3	11
29	2	12	3	11
30	2	12	3	11