



Serie documentos de trabajo

EL USO DE MODELOS LOG-LINEALES PARA EL ANÁLISIS DEL CONSUMO RESIDENCIAL DE ENERGÍA

Enrique de Alba ITAM

Yolanda Mendoza Dirección General de Política Energética SEMIP

DOCUMENTO DE TRABAJO

Núm. VI - 1984

EL USO DE MODELOS LOG-LINEALES PARA EL ANALISIS DEL CONSUMO RESIDENCIAL DE ENERGIA

Enrique de Alba

У

Yolanda Mendoza

Este trabajo está dirigido a quienes se interesen en la aplicación de métodos estadísticos en Economía. Se presenta primero - una reseña breve de lo que es el modelo log-lineal, su defini -- ción, algunas características, métodos de estimación, selección- de modelos, medidas de asociación y análisis de los resultados.- Al final se aplica la metodología a un cuadro de resultados de - la Encuesta de Ingreso y Gasto de los Hogares, 1977 de la Secretaría de Programación y Presupuesto, para ilustrar el uso del modelo. Los resultados sirven para verificar algunas hipótesis respecto al comportamiento de las familias en cuanto al gasto en -- energéticos.

ANALYSING RESIDENTIAL ENERGY CONSUMPTION AN APPLICATION OF LOG-LINEAR MODELS

Enrique de Alba

and

Yolanda Mendoza

This article is intended for those who are interested in applying statistical methods in Economics. A brief presentation is made of the log-linear model, its definition, some properties, estimation methods, model selection, measures of association between variables and interpretation of the results. In the last part an application is made to a table of results from the 1977 Household Income and Expenditure Survey, carried out by Secretaría de Programación y Presupuesto, to illustrate the use of this kind of model. The results obtained verify some hipotheses on household behavior regarding expenditures on energy.

EL USO DE MODELOS LOG-LINEALES PARA EL ANALISIS DEL CONSUMO RESIDENCIAL DE ENERGIA

ENRIQUE DE ALBA

Instituto Tecnológico Autónomo de México

Y

YOLANDA MENDOZA

Dirección General de Política Energética SEMIP

INTRODUCCION

Los métodos econométricos que se utilizan con mayor frecuencia, o por lo menos los que ocupan un mayor espacio en los textos sobre el tema, se basan en el supuesto de variables dependientes contínuas; así tenemos el caso del modelo lineal y series de tiempo. En ocasiones pueden presentarse datos correspondientes a variables cua litativas, como la distinción entre zonas urbanas ó zonas rurales; también puede suceder que a pesar de que las variables sean continuas los resultados se presenten agrupados, como por ejemplo estratos de ingreso familiar. En estas situaciones los métodos anteriores no son adecuados y habrá que recurrir a métodos para el análisis de datos categóricos.

En esta nota se pretende hacer una exposición sencilla de los conceptos que se utilizan para el análisis estadístico de datos categóricos, concretamente mediante el uso del modelo conocido como log -lineal.

Se ha dicho que las variables se dividen en categorías.

Las categorías se presentan cuando clasificamos a la población, digamos en Q grupos. Se dice una clasificación es exhaustiva cuando pue den clasificarse a todos los miembros de la población en alguna categoría; dichas categorías son mutuamente exclusivas cuando están definidas de tal manera que se pueda acomodar a un indi---

viduo correctamente en una y sólo una de ellas. Cuando se hace una clasificación se busca que tenga estas dos propiedades, para que no haya confusión al acomodar a un individuo de población y caiga siem pre en algún grupo.

Al hablar de datos cualitativos es necesario un arreglo que refleje la estructura de los datos. Este arreglo debe definir las categorías de cada variable y las interrelaciones con las demás categorías de las otras variables en el modelo. Un arreglo rectangular es la estructura más adecuada para representar a 2 variables: los renglones del arreglo se hacen corresponder con los niveles de la la. variable mientras que las columnas se asocian con los niveles de la 2a. variable. Si se desea presentar más de 2 variables pueden utilizarse tantos arreglos rectangulares como sea necesario. Por ejemplo, supóngase que se tienen como variables el ingreso (alto, medio, bajo) y el sexo (masculino, femenino); se presenta la información en una tabla, los renglones correspondientes a cada nivel de ingreso y las columnas al sexo:

		SEXO Masculino Femenino
	Alto	
Ingreso	Medio	
	Bajo	
	Dajo	

Si se considera la inclusión de una tercera variable por ejem plo el tipo de población en el que viven los individuos (rural, urbana) la información se resumiría en 2 arreglos idénticos al anterior, uno para la población rural y otro para la población urbana.

		Població	on Rural	Población Urbana		
		Se	хо		Se	хо
		Masculino	Femenino		Masculino	Femenino
Ingreso	Alto			Alto		
	Medio			Medio		·
	Bajo		: :	E Bajo		

TABLA 2

Al observar el arreglo vemos que la posición en las celdas nos refleja las características de los individuos que caen en ellas, a este arreglo se le conoce como tabla de contingencia, en datos cualitativos multivariados cada individuo se describe por una serie de atributos; al tomar una mues tra todos los individuos con la misma descripción son enumerados y esta cuenta entra en la celda correspondiente de la tabla de contingencia, por ejemplo, todos los individuos en la muestra que tengan ingreso alto y sean hombres entran en la celda superior izquierda de la tabla 1.

Una vez que se tiene la tabla de contingencia el interés se enfoca

en analizar las relaciones que existen entre las variables. El modelo log-lineal proporciona un esquema para el análisis de dichas tablas.

En el análisis de tablas de contingencia se presentan dos situaciones:

- i) Una variable se vé como de respuesta y las demás son explicativas, esto es lo que se conoce como tabla asimétrica.
- ii) No se hace distinción entre variables dependientes e independientes, esto es lo que se llama tabla simétrica.

El objetivo del Análisis es plantear los dos tipos de situaciones, primero no distinguiendo las variables, o sea considerándolas conjuntamente y posteriormente planteando un modelo en el que se planteen relaciones de causalidad.

Además, para el análisis de las relaciones entre 2 variables se utiliza una medida de asociación, que se presenta para cada tabla bivariada y nos dá una idea de como están asociadas las variables de dicha tabla, la magnitud de dicha asociación, así como si es o no significativa.

I. Modelo bivariado

I.A) El caso 2 x 2

Se ha planteado la importancia del análisis de datos categóricos y la idea de la exposición es plantear un modelo con un número cual lesquiera de variables, sin embargo, antes de generalizar el modelo, se expondrá la tabla de contingencia más simple, la tabla 2 x 2, basada en 4 celdas, con dos categorías para cada variable, pues esto facilitará la comprensión de un modelo con un mayor número de varia bles

Considérese la variable A_1 que puede clasificarse en las categorías i_1 (i_1 =1,2) y la variable A_2 que se clasifica en las categorías i_2 (i_2 =1,2). Supongamos que los dos renglones de la tabla corresponden a las categorías de la variable 1 (A_1) y que las columnas se asocian a la variable 2 (A_2). Se utilizan subíndices para referirse a la posición en la tabla, el 10. le corresponde a la variable 1 y el 20. a la variable 2.

El modelo está basado en probabilidades (p) pero para referirse a frecuencias se utiliza una transformación simple, como se verá más adelante. Considérese la tabla:

		1	2 2
$\mathbf{A_1}$	1	P_{11}	P ₁₂
1	2	P ₂₁	P ₂₂

En donde $\mathbf{p}_{i_1 i_2}$ es la probabilidad que tiene un individuo de estar en la categoría i_1 de la variable 1 y en la categoría i_2 de la variable 2.

El modelo log-lineal para esta tabla se escribe :

$$\log p_{i_1 i_2} = \mu + \alpha_1(i_1) + \alpha_2(i_2) + \beta_{12}(i_1, i_2)$$

$$i_1 = 1, 2$$

$$i_2 = 1, 2$$

a μ , α_1 (i_1), α_2 (i_2) y β_{12} (i_1 , i_2) se les conoce como parámetros o efectos del modelo. μ es la media general de los logarít mos de las probabilidades, es decir :

$$\mu = \frac{1}{4} \left(\log p_{11} + \log p_{12} + \log p_{21} + \log p_{22} \right)$$

 $\mu + \alpha_1$ (i_1) es la media de los logarítmos de las probabili-

dades en el nivel i_1 de la variable 1 y se obtiene de :

$$\mu + \alpha_1(i_1) = \frac{1}{2} (\log p_{i_1} + \log p_{i_1})$$

$$i_1 = 1, 2$$

 $\mu + \alpha_2$ (i_2) es la media de los logarítmos de las probabilidades en el nivel i_2 de la variable 2.

$$\mu + \alpha_2(i_2) = \frac{1}{2} [log p_{1i_2} + log p_{2i_2}]$$

$$i_2 = 1, 2$$

Notese que $\beta_{12}(i_1,i_2)$ puede encontrarse mediante simples operaciones algebraicas una vez que se conocen los valores de μ , $\alpha_1(i_1)$ y $\alpha_2(i_2)$, pues

$$\beta_{12}(i_1, i_2) = \log p_{i_1i_2} - \alpha_1(i_1) - \alpha_2(i_2)$$

$$= \log p_{i_1i_2} - \mu - \frac{1}{2} \sum_{i_1} \log p_{i_1i_2}$$

 $+ \mu - \frac{1}{2} \sum_{i_2} \log p_{i_1 i_2} + \mu$

sustituyendo μ por su valor como media general, se obtiene :

$$\beta_{12}(i_1, i_2) = \log p_{i_1 i_2} - \frac{1}{2} \sum_{i_1} \log p_{i_1 i_2}$$

$$- \frac{1}{2} \sum_{i_2} \log p_{i_1 i_2} + \frac{1}{4} \sum_{i_1 i_2} \log p_{i_1 i_2}$$

A continuación se presentan las restricciones a que estan sujetos estos parámetros.

I.A.1) Restricciones ANOVA

Como $\alpha_1(i_1)$ y $\alpha_2(i_2)$ representan desviaciones con respecto a la media general, al sumar sobre todos los valores de un subíndice de éstos efectos la suma es cero, lo cual se expresa :

$$\alpha_1(1) + \alpha_1(2) = \alpha_2(1) + \alpha_2(2) = 0$$

Si se define $\ell_{i_1i_2} = \log p_{i_1i_2}$ se escriben los parámetros como :

$$\mu = \frac{1}{4} \sum_{i_1 i_2} \ell_{i_1 i_2} = \frac{1}{4} \ell_{\bullet \bullet}$$

$$\alpha_1(i_1) = \frac{1}{2} \ell_{i_1 \bullet} - \frac{1}{4} \ell_{\bullet \bullet}$$

$$\alpha_2(i_2) = \frac{1}{2} \ell_{\bullet i_2} - \frac{1}{4} \ell_{\bullet \bullet}$$

en donde el punto indica suma sobre los subíndices donde aparece Esta notación permite que se muestre que se cumplen las restricciones, véase que al sumar sobre i_1 en α_1 se obtiene:

$$\hat{\alpha}_1(\cdot) = \sum_{i_1=1}^2 \left(\frac{\ell_{i_1}}{2} - \frac{\ell_{\cdot}}{4} \right) = 0$$

de la misma manera se cumple con que: $\alpha_2(\cdot) = 0$

El término β_{12} representa una desviación con respecta a $\mu + \alpha_1(i_1) + \alpha_2(i_2)$ o sea:

$$\beta_{12}(i_1, i_2) = \ell_{i_1 i_2} - \frac{\ell_{i_1 \bullet}}{2} - \frac{\ell_{\bullet i_2}}{2} + \frac{\ell_{\bullet \bullet}}{4}$$

así que,

$$\beta_{12}(i_1, \cdot) = \beta_{12}(\cdot, i_2) = 0$$
; $i_1 = 1, 2$
 $i_2 = 1, 2$

o dicho de otra manera

$$\beta_{12}(1,1) = -\beta_{12}(1,2) = -\beta_{12}(2,1) = \beta_{12}(2,2)$$

I.A.2) Modelos con frecuencias celdales

Los modelos pueden describir, como se dijo, frecuencias celdales esperadas en lugar de probabilidades; la diferencia en la formulación estriba en una constante. Si se considera una tabla con conteos esperados m , de tal manera que el tamaño de muestra es :

$$n = \sum_{i_1 i_2} m_{i_1 i_2}$$

donde

$$m_{i_1i_2} = n p_{i_1i_2} ;$$

tomando logarítmo,

$$\log m_{i_1 i_2} = \log n + \ell_{i_1 i_2}$$

y entonces el modelo log-lineal considerando frecuencias esperadas celdales es ($log m_{i_1i_2} = \ell_{m_{i_1i_2}}$):

$$\ell_{m_{i_1 i_2}} = \mu' + \alpha_1(i_1) + \alpha_2(i_2) + \beta_{12}(i_1, i_2)$$

con

$$\mu' = \log n + \mu$$

Se utilizara μ en cualquiera de los casos, ya sea que se trate de frecuencias esperadas o probabilidades. En la siguiente sección se plantea la construcción del modelo a partir de la condición de independencia.

I.A.3) Formulación del modelo a partir de la condición de - independencia

La tabla de contingencia para las variables 1 y 2, puede extenderse a incluir los totales marginales, los cuales se expresan:

$$\sum_{\mathbf{i_2}} \qquad \mathbf{p_{i_1 i_2}} = \mathbf{p_{i_1}}$$

$$\sum_{i_1} - p_{i_1 i_2} = p_{i_2} ,$$

que son los totales marginales renglón y columna respectivamente, dan la probabilidad que tiene un individuo de caer en la categoría i_1 de la variable 1, y la probabilidad que tiene de caer en la categoría i_2 de la variable 2, respectivamente.

El arreglo completo es:

		1	2	· :	
Α.	1	P _{1 1}	P _{1 2}	P ₁ .	
A ₁	2	. P _{2 1}	P _{2 2}	P ₂ .	
		P. 1	P. ₂	1	

Para llegar a la formulación del modelo a la manera log-lineal se supone que A_1 y A_2 son independientes. La condición de independencia se expresa:

$$p_{i_1 i_2} = p(A_1 = i_1) * p(A_2 = i_2) = p_{i_1 \cdot p_{\bullet i_2}}$$

Utilizando conteos esperados,

$$m_{i_1 i_2} = \frac{m_{i_{10}}}{n}$$

Del modelo log-lineal en donde se ha omitido el parámetro β_{12} (i_1 , i_2) se sabe que;

$$m_{i_1} = e^{\mu + \alpha_2(i_2)} \sum_{i_1} e^{\alpha_1(i_1)}$$

$$m_{i_2} = e$$

$$\sum_{i_2}^{\mu} \alpha_2(i_2)$$

$$n = e^{\mu} \sum_{i_1} e^{\alpha_1(i_1)} \sum_{i_2} e^{\alpha_2(i_2)}$$

de donde:

$$\mu + \alpha_{1}(i_{1}) + \alpha_{2}(i_{2})$$

$$= m_{i_{1}i_{2}}$$

Por lo tanto, el modelo log-lineal que supone independencia es consistente con la formulación que se hace en tablas de contingencia de esta condición.

I.B) Extensión a tablas I₁ x I₂

Se puede extender la tabla para tener I_1 categorías para la variable l e I_2 categorías para la variable l. Los parámetros en este contexto se calculan como sigue, la medida general es:

$$\mu = \frac{1}{I_1 I_2} \sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} \ell_{m_{i_1} i_2}$$

los efectos principales se escriben

$$\alpha_1(i_1) = \frac{1}{I_2} \sum_{i_2=1}^{I_2} \ell_{m_{i_1 i_2}} - \mu$$

$$\alpha_2(i_2) = \frac{1}{I_1} \sum_{i_1=1}^{I_1} \ell_{m_{i_1}i_2} - \mu$$

y el efecto bivariado se expresa:

$$\beta_{12}(i_1,i_2) = \ell_{m_{i_1i_2}} - \frac{1}{l_2} \sum_{i_2=1}^{l_2} \ell_{m_{i_1i_2}} - \frac{1}{l_1} \sum_{i_1=1}^{l_1} \ell_{m_{i_1i_2}} + \mu$$

I.C) Interpretación de parámetros

Los parámetros del modelo log-lineal son análogos a los efectos en Análisis de Varianza (ANOVA) cuando una variable se descompone en términos aditivos que representan una media general, efectos principales y efectos de orden mayor (bivariados, trivariados, etc.)

El término μ es un efecto normalizante que permite que:

$$\sum_{i_1=1}^{l_1} \sum_{i_2=1}^{l_2} p_{i_1i_2} = 1$$

o bien:

$$\sum_{i_1=1}^{I_1} \sum_{i_2=1}^{I_2} m_{i_1 i_2} = r$$

El término $\alpha_i(i_i)$ representa el efecto principal de la log-frecuencia de estar en el nivel i_i de la variable l.

Similarmente $\alpha_2(i_2)$ representa el efecto principal de la log-frecuencia de estar en la categoría i_2 de la variable 2.

El término $\beta_{12}(i_1,i_2)$ representa el efecto conjunto de la log-frecuencia de estar en la categoría i_1 de la variable 1 y en la categoría i_2 de la variable 2.

La generalización a tablas de contingencia con más de dos variables sigue un desarrollo análogo al presentado, con la utilización de algebra matricial para facilitar la exposición.

II. Modelo Log-Lineal q-variado.

La formulación del modelo log-lineal se ha presentado como una descomposición de las log-frecuencias de una tabla de contingencia en varios componentes aditivos. Siguiendo esta idea se presenta a continuación la generalización de este modelo a una tabla multidimensional que incluya q variables.

Supóngase que se tienen A_1,A_2,\ldots,A_q variables categóricas, una tabla de contingencia consiste de $x_{i_1\,i_2\,\ldots\,i_q}$ frecuencias donde

$$i_{j} = 1, 2, \dots, I_{j}$$
, con $j = 1, 2, \dots, q$

donde I_j es el total de categorías para la variable j. Se puede simplificar a un sólo subindice i, que se forma de la siguiente manera:

$$i = (i_1 - 1) I_2 I_3 ... I_q + (i_2 - 1) I_3 I_4 ... I_q + ...$$

$$+ (i_{q-1} - 1) I_q + i_q ...$$

El uso de este subíndice permite asignar un sólo número a cada una de las celdas, sin embargo, no nos indica las caracterís ticas de los individuos que caen en ellas, es un índice lexicográfi-

co, es decir, sigue un cierto orden (2). Se puede escribir entonces:

$$x_{i} = x_{i_{1} i_{2} \cdots i_{q}}$$
, $i = 1, 2, \ldots, Q$

tal que:

$$Q = \prod_{k=1}^{q} I_k$$

es el número total de celdas que tiene la tabla.

Para facilitar la explicación de algunos conceptos que se plantearán más adelante, es conveniente poder referirse a un conjunto cualquiera de subíndices:

$$\eta = \left\{ i_1, i_2, \ldots, i_q \right\}$$

en donde los subíndices toma los valores

$$i_1 = 1, 2, ..., I_1$$
 $i_2 = 1, 2, ..., I_2$
 $\vdots \vdots \vdots$
 $i_q = 1, 2, ..., I_q$

$$\sum_{n} p_{n} = 1$$

Para poder definir a los modelos log-lineales en un contexto general, se requiere su expresión en algebra matricial, por lo cual se ordena los log p_{-} s en un vector:

$$\underline{\ell} = \underline{\log p} = \begin{bmatrix}
\ell_{11} & \dots & 1 \\
\ell_{11} & \dots & 2
\end{bmatrix}$$

$$\ell_{1_1 1_2 \dots 1_q}$$

En donde las ℓ 's se particionan a la manera ANOVA como:

$$\ell_{\eta} = \log p_{\eta} = \mu + \alpha_{1}(i_{1}) + \ldots + \alpha_{q}(i_{q}) + \beta_{12}(i_{1}, i_{2}) + \ldots + \beta_{q-1}(i_{q-1}, i_{q}) + \ldots + \beta_{q-1}(i_{q-1}, i_{q}) + \ldots$$

 $+ \omega_{12...}$ (i_1, i_2, \ldots, i_g)

o bien:

$$\log p_{\eta} = \mu + z_{\eta}$$

donde μ denota un efecto global; $\alpha_1(i_1)$ es el efecto debido a A_1 en el nivel i_1 , se le llama efecto principal de A_1 ; $\alpha_q(i_q)$ es el efecto debido a A_q en el nivel i_q ; $\beta_{12}(i_1,i_2)$ denota la interacción de 2^2 orden o efecto bivariado entre A_1 y A_2 a los niveles i_1 , i_2 respectivamente y así succesivamente hasta que $\omega_{12}\ldots_q$ (i_1,i_2,\ldots,i_q) es el efecto de interacción de orden q o efecto q - variado entre las q variables A_1 , A_2,\ldots,A_q a los niveles i_1 , i_q .

Como $p_{\eta} < 1$, log P_{η} es negativo, pero μ no es fijo y debido a esto los efectos no están restringidos en signo.

Las restricciones ANOVA exigen que si se suman los efectos sobre todos los valores de un subíndice, la suma debe ser cero, por ejemplo:

$$\alpha_{1}(\cdot) = \sum_{i_{1}=1}^{l_{1}} \alpha_{1}(i_{1}) = 0$$

El punto denota suma sobre el subíndice donde aparece.

Se utiliza η para denotar un conjunto que incluya a las q variables, y η , un subconjunto que incluya a ciertas variables, . Se puede indexar una celes un número y hablar de la frecuencia observada en la celda . Por ejemplo supóngase que el conjunto de variables $\sigma = \{A_1, A_2, ..., A_q\}$ se particiona en 2 subconde tal manera que σ_1 , σ_2 juntos $\sigma_1 = \{A_1, A_2, \ldots, A_s\}$, $\sigma_2 = \{A_{s+1}, A_{s+2}, \ldots, A_q\}$ η_1 y η_2 para denotar subcon-Entonces puede utilizarse juntos de sodices que se refieran a las categorías de las variables respectivamente. incluídas en σ_1 σ_2 y

Dadas las convenciones anteriores, la probabilidad de caer en una celda cualquiera se expresa:

$$p_{\eta} = p_{i_1 i_2 \dots i_q} = p_r \left(A_1 = i_1, A_2 = i_2, \dots, A_q = i_q \right)$$

y representa la probabilidad de caer en la celda η . La formulación del modelo permite que se cumpla con que la suma de las probabilidades sea l, o sea:

Se requiere que se cumplan el siguiente conjunto de restric ciones:

$$\alpha_1(\cdot) = \ldots = \alpha_q(\cdot) = 0$$

$$\alpha_{1}(\cdot) = \dots = \alpha_{q}(\cdot) = 0$$

$$\beta_{12}(\cdot, i_{2}) = \dots = \beta_{q-1} \cdot (i_{q-1}, \cdot) = 0$$

$$\omega_{12...q}(\cdot,...,i_q) = ... = \omega_{12...q}(i_1,...,\cdot) = 0$$

Un ejemplo de la expresión de la matriz de diseño en el caso saturado 2 x 2 para las variables A, y A2

$$\begin{cases} \mu \\ \alpha_{1}(1) \\ \alpha_{1}(2) \\ \alpha_{2}(1) \\ \alpha_{2}(2) \\ \beta_{12}(1,1) \\ \beta_{12}(1,2) \\ \beta_{12}(2,1) \\ \beta_{12}(2,2) \end{cases}$$

En esta expresión no se han incluído las restricciones tipo ANOVA por lo que se observa que la matriz de diseño, la que contiene ceros y unos, es singular. Al incorporar las restricciones el modelo queda de la siguiente forma:

III. Estimación

Introduccion.

res ajustados para cada una de las celdas en la tabla de contingen cia, así como los valores de los parámetros del modelo log-lineal. Para calcular estos valores se supone cierto comportamiento teórico caracterizado por el valor de la probabilidad de caer en las celdas; este valor depende de los parámetros incluídos dentro del modelo.

La tabla completa consta de celdas elementales y en ella, puede sumarse sobre uno o varios subíndices, estas sumas se \underline{agru} pan en tablas de celdas no elementales que tienen menos dimensiones que el arreglo completo. Llamamos a estas tablas de sumas configuraciones y se denotan mediante C. Por ejemplo, supóngase que se tienen tres variables, el arreglo bidimensional obtenido de sumar sobre la 3a. variable, el conjunto $\{x_{i_1i_2\bullet}\}$ es la configuración C_{12} . Cada una de las $i_1 \times i_2$ celdas no elementales está formada por la suma de i_3 celdas elementales, dando como resultado que la configuración C_{12} involucra unicamentales

te a las variables l y 2.

Para ilustrar lo anterior, supóngase que se tienen datos acerca del envenenamiento después de un día de campo. Las variables son:

- 1) Envenenamiento { Presente No presente
- 2) Ensalada de papas consumida { si no
- 3) Jaiba consumida { si no

La muestra consiste de 304 personas que asistieron àl día de campo y los resultados del cuestionario que llenaron se presentan en la siguiente tabla: (6)

Jaiba

Comida consumida

Jaiba

	si		no •	
	Ensalada de papas		Ensalada	de papas
Envenenamiento	si	no	si	no
si	120	4	20	0
no	80	31	24	23

Esta tabla consta de celdas elementales; para interpretarla puede observarse que según esta tabla 4 personas de la muestra comieron jaiba, no comieron ensalada y presentaron envenenamiento. La configuración C_{12} representa la suma sobre la 3a. variable, jaiba, y es:

	Ensalada de papas		
Envenenamiento	si	no	
si	142	4	
no	104	54	

Nótese que ha desaparecido la variable jaiba, cada una de las celdas representa la suma sobre los que comieron y no comieron jaiba. Por ejemplo, el valor 142 son los individuos que comieron ensalada y presentaron envenemento o sea 120 + 22.

Otro concepto importante dentro del análisis de tablas de contingencia es el de tablas incompletas, que se definen en términos del tipo de celdas vacias (con valor cero) que tengan. Los ceros de una tabla pueden clasificarse en:

- i) Ceros muestrales o aleatorios: Surgen cuando la mues tra no es lo suficientemente grande para obtener respuesta en una combinación particular de características.
- ii) Ceros fijos o estructurales: Cuando una combinación específica de atributos no ocurre en una población, estas celdas deben permane-

cer vacias, y esto a menudo se sabe a priori.

Las tablas con ceros estructurales se llaman incompletas y su análisis es complicado ya que debe imponerse la restricción que ciertas probabilidades sean cero.

Se limitará el estudio a las tablas completas que pueden o no presentar ceros muestrales. La presencia de ceros muestrales puede provocar que un modelo planteado no pueda estimarse; existen reglas para detectar cuando un modelo no es estimable (1).

Por otro lado, al escoger un modelo debe tenerse en cuenta que hay 2 límites del rango de modelos a ajustar.

Estos limites son:

- i) El modelo saturado: en donde las log-frecuencias ajustadas encajan exactamente en las log-frecuencias observadas. Se tienen cero grados de libertad y la estadística que mide la bondad del ajuste toma el valor de 0.
- ii) El modelo mínimo: contiene el conjunto mínimo de parámetros permitido. Para ciertos esquemas muestrales
 (Poisson y multinomial) solo incluye a μ, mien---

⁽¹⁾ Estas reglas están basadas en el examen de algunas configuraciones que surgen a partir de la tabla completa. Véase Kawasaki (1979) Capítulo IV.

tras que para otros (multinomial condicional y otros) incluye solamente a los parámetros asociados con las marginales fijas por el modelo.

El objetivo del análisis es encontrar un buen modelo situado entre éstos extremos que nos describa las relaciones que existen en tre las variables del modelo. Es deseable que este modelo posea parámetros estadísticamente significativos, que presente un buen ajuste y que sea simple, ésto último se refiere a incluir el menor número de parámetros posible.

III.A Método de Estimación de Máxima Verosimilitud.

En esta sección se plantea la estimación del modelo loglineal mediante el Método de Máxima Verosimilitud. Este implica que debe haber algún supuesto acerca de la distribución que genera las ob servaciones (frecuencias) en las celdas, es decir se supone una forma específica para la función de densidad $f(x_i)$ donde x_i es la frecuencia en la i - ésima celda y f es una función de densidad -- - (f.d.p.).

Los esquemas muestrales que más frecuentemente se utilizan son:

- a) Muestreo Poisson independiente
- b) Muestreo Multinomial simple
- c) Muestreo Multinomial condicional

El primero, Poisson, surge cuando para un período fijo se observa un conjunto de procesos Poisson y no está fijo el tamaño de la muestra. El multinomial se presenta cuando se fija el tamaño de muestra. El último surge cuando se imponen condiciones a los to tales marginales.

Una vez planteado el modelo, se deberá maximizar la función Verosimilitud lo cual se logra mediante algún método numérico. A continuación se presenta el de ajuste proporcional iterativo.

III.B) Método de Ajuste Proporcional Iterativo.

Para aplicar este método es necesario establecer las estadísticas suficientes que corresponden al modelo log-lineal que va a ajustarse, se limita el ajuste a modelos log-lineales jerárquicos, es decir, modelos en los que la presencia de un término de orden mayor implica que los términos asociados a este de menor orden esten presentes, o bién, si un término de menor orden no se incluye, entonces todos los correspondientes términos de mayor orden relacionados con él están ausentes.

Las estadísticas suficientes son las configuraciones asociadas a cada parámetro del modelo log-lineal. Por ejemplo en el caso de 2 variables la configuración asociada al parámetro β_{12} es C_{12} .

Como se trata de modelos log-lineales jerárquicos es posible establecer un conjunto mínimo de estadísticas suficientes en donde se incluyen todas aquellas configuraciones que correspondan a los términos de orden mayor (s) en el modelo y examinar si a partir de este conjunto es posible obtener las demás estadísticas suficientes asociadas a los parámetros.

Si no, se examinan los términos de orden s - 1 que no estén

asociados a los de orden s, y se incluyen en el conjunto las configuraciones asociadas a estos.

Se continúa así hasta tener en el conjunto mínimo de estadísticas suficientes las configuraciones necesarias, para que a partir de éstas se obtengan todas las estadísticas suficientes asociadas con los parámetros del modelo log-lineal.

Por ejemplo, a partir de la configuración C_{12} es posible obtener la configuración C_1 . y $C_{\cdot 2}$.

El procedimiento que se describe a continuación parte de que, una vez definido el conjunto de estadísticas mínimas suficientes, las frecuencias esperadas del modelo se calculan igualando los totales marginales esperados con los valores correspondientes de las estadísticas suficientes.

Aunque en ocasiones puede ocurrir que las estimaciones en las celdas se escriben como una función directa de las estadísticas suficien tes, lo que se denomina estimación directa; si se trabaja con un conjun to razonablemente grande de datos, no se requiere determinar si existe o no estimación directa, ya que este método da las estimaciones automáticamente. El procedimiento iterativo que se expone ajusta las configuraciones suficientes ciclo a ciclo y posee las siguientes propiedades:

- i) Siempre converge al conjunto único de estimaciones de M-V.
- ii) Permite especificar el grado de precisión de las estimaciones en las celdas elementales, además de que éstos valores dependen de las estadísticas suficientes.
- iii) Se puede escoger cualquier conjunto de estimaciones iniciales y si hay estimaciones directas el procedimiento converge en un ciclo.

En general, mediante éste método las celdas internas de la tabla son proporcionalmente ajustadas a un conjunto de marginales. La convergencia queda asegurada ya que la verosimilitud es una función monótona decreciente por lo que siempre pueden obtenerse estimaciones con cierta precisión fija de antemano.

La presencia de celdas vacías causa que todas aquellas celdas que formen una configuración mínima con alguna celda vacía, tendrán como estimación cero. Supongáse por ejemplo un modelo con tres variables en donde se va a ajustar la configuración mínima C_{12} y $x_{11} = 0$, entonces los valores ajustados para cada una de las celdas elementales que sumadas son x_{11} serán cero, es decir m_{11} i_3 i_3 i_3 i_4 i_5 i_6 i_6 i_7 i_8 i_8 i_8 i_8 i_8 i_9 i_9

1) Descripción del método

Supóngase que se han escogido como estadísticas mínimas

suficientes las configuraciones: C_{η_i} , donde t=1,..., s, con conteos celdales marginales x_{η_i} respectivamente. Se escoge un conjunto de estimaciones iniciales $m_{\eta}^{(0)}$ y se ajusta cada configuración en turno, utilizando en cada paso del ciclo la estimación obtenida en el paso anterior. Como hay que ajustar cada una de las configuraciones en cada ciclo, en el primer ciclo se ajustan s configuraciones, en el segundo ciclo s configuraciones más, dando un total de 2 s configuraciones ajustadas. Por lo tanto las configuraciones que se han ajustado al término de r ciclos son rs. Es por esto que al iniciar un nuevo ciclo las relaciones son:

$$\widehat{\mathbf{m}}_{\eta}^{(rs+1)} = \widehat{\mathbf{m}}_{\eta}^{(rs)} \frac{\mathbf{x}_{\eta_1}}{\widehat{\mathbf{m}}_{\eta_1}^{(rs)}}$$

Siendo C_{η_1} la primera configuración a ajustar, o sea, el procedimiento vuelve a ajustar las s configuraciones consecutivamente. Al comienzo de este ciclo se tiene la estimación recien obtenida que es la que se utiliza para obtener una nueva estimación. La estimación correspondiente será la rs+1 y se vuelven a ajustar todas las configuraciones, empezando con la η_1 . En general el K - ésimo es:

$$\widehat{\mathbf{m}}^{(K)} = \widehat{\mathbf{m}}^{(K-1)} \frac{\mathbf{x}\eta_t}{\widehat{\mathbf{m}}^{(K-1)}}$$

se sugiere utilizar como estimaciones iniciales $\widehat{\mathbf{m}}^{(0)} = 1$

2) Convergencia del método

Para estipular en que momento parar, es necesario considerar el cociente de verosimilitud después de ajustar k configuraciones:

$$D_{\eta}^{(K)} = \sum_{x \in \eta} \log_{x \in \eta} - \sum_{x \in \eta} \log_{\eta} \widehat{m}_{\eta}^{(K)}$$

lo que interesa es maximizar la verosimilitud y la prueba de convergencia depende del cambio que ocurra en ésta. Las estimaciones ma
ximizan el kernel de la log-verosimilitud o sea Σ x η log $\widehat{\mathbf{m}}_{\eta}$

sujeto a restricciones, éste valor aumenta a medida que aumenta k, y la frontera superior es Σ x η log x η ; por lo tanto, un buen criterio de convergencia es $D_{\eta}^{(K)}$. Debe determinarse que tán pequeña debe ser la diferencia entre ciclos para ser considerada insignificante. La regla para parar es estipular que los cambios ciclo a ciclo no excedan un valor pequeño δ .

Para el caso considerado de s pasos, es decir, s configuraciones mínimas, el cambio en la función de verosimilitud entre ciclos es:

$$\triangle \mathbf{p} = \mathbf{p}_{\eta}^{(K-S)} - \mathbf{p}_{\eta}^{(K)}$$

IV. Diversas Pruebas

Cuando se ajusta un modelo normalmente se quiere tener alguna idea de la bondad de ajuste del mismo, la confiabilidad de los parámetros estimados, grado de asociación entre variables, etc. En el caso de los modelos log-lineales también son de interés algunos indicadores relativos a la confiabilidad de las estimaciones celdales.

A continuación se presenta algunos criterios para llevar a cabo lo anterior en el caso de los modelos log-lineales.

IV.A) Estadísticas de bondad de ajuste.

Para conocer la significancia estadística de un modelo, es decir, qué tanto se aproximan las estimaciones a los datos observados, se requiere obtener ciertas estadísticas. Estas estadísticas nos describen que tan bueno es un modelo, o sea su bondad de ajuste, comparando las frecuencias observadas contra las esperadas bajo un esquema particular.

Las estadísticas χ^2 y G^2 que se describen a continuación tienen una distribución χ^2 asintótica cuando el modelo es correcto. Los grados de libertad de estas estadísticas pueden obtenerse de 2 maneras:

- i) Contando el número total de parámetros independientes que se hacen cero.
- ii) Contando el número total de parámetros estimados y restando este número del total de celdas estimadas.

I.A) ___ Ji - Cuadrada de Pearson.

Se define como:

$$X^{2} = \sum \frac{(x_{i} - \widehat{m}_{i})^{2}}{\widehat{m}_{i}^{2}}$$

donde:

x; es la observación en la celda i - ésima

 $\hat{\mathbf{m}}_{i}$ es la estimación para la celda i - ésima.

Esta estadística es una medida del ajuste global del modelo y es muy utilizada en el estudio de tablas de contingencia.

IV.B) Cociente de Verosimilitud.

Se expresa como:

$$G^2 = -2 \sum_{i=1}^{\infty} x_i \log_{i} \left(\frac{\widehat{m}_i}{x_i} \right)$$

que es menos dos veces el cociente de verosimilitud. El valor de G^2 se interpreta como la probabilidad de que las diferencias en tre las frecuencias observadas y ajustadas hayan sido aleatorias dado que el modelo sea correcto.

La estadística G^2 tiene una propiedad interesante que es la de poder descomponerse condicionalmente. Esto es, supóngase que se tiene un modelo (2) anidado dentro de otro modelo (1), es decir, el modelo 2 contiene una parte de los parámetros del modelo 1. Entonces la G^2 del modelo 2, denotada G^2 (2) que mide el ajuste de las estimaciones bajo el modelo 2,

puede descomponerse en:

- i) Una medida del ajuste del modelo 1,
- ii) Una medida condicional del modelo 2 dado el modelo l. Dicho de otra manera, una medida de las distancias de las estimaciones $\{\widehat{m}_i^{(2)}\}$ a las estimaciones $\{\widehat{m}_i^{(1)}\}$ que se denota G^2 (2/1).

Para ilustrar lo anterior, se expone a continuación la descomposición de $\mbox{ G}^{2}$ (2):

$$G^{2}(2) = -2 \sum_{i} x_{i} \qquad \left(\frac{\widehat{m}_{i}^{(2)}}{x_{i}}\right)$$

$$= -2 \sum_{i} x_{i} \log \widehat{m}_{i}^{(2)} + 2 \sum_{i} x_{i} \log \widehat{m}_{i}^{(1)}$$

$$-2 \sum_{i} x_{i} \log \widehat{m}_{i}^{(1)} + 2 \sum_{i} x_{i} \log x_{i}$$

$$= -2 \sum_{i} x_{i} \log \left(\frac{\widehat{m}_{i}^{(2)}}{\widehat{m}_{i}^{(1)}}\right) - 2 \sum_{i} x_{i} \log \left(\frac{\widehat{m}_{i}^{(1)}}{x_{i}}\right)$$

$$= G^{2}(2/1) + G^{2}(1)$$

en esta expresión se observan las partes ii) e i) mencionadas que forman \mathbf{a} $\mathbf{G}^2(2)$.

De este resultado puede obtenerse:

$$G^{2}(2/1) = G^{2}(2) - G^{2}(1)$$

que puede utilizarse para verificar si los parâmetros del modelo 2 le añaden significancia al modelo 1. Mediante esta fórmula pueden compararse los 2 modelos y lo que se prueba es si la adición de parâmetros resulta en una reducción significativa de la G^2 (2) . Para este efecto se calcula G^2 (2) - G^2 (1) que se distribuye como una χ^2 con $\nu_2 - \nu_1$ grados de libertad, si se supone que los grados de libertad respectivos de G^2 (2) y G^2 (1) son ν_2 y ν_1 .

IV.C) Transformación de Freeman - Tukey.

Cuando la muestra es pequeña no puede decirse cómo se distribuyen las estadísticas anteriores, ya que X^2 y G^2 se distribuyen como una X^2 asintóticamente, es decir se aproximan a una X^2 cuando el tamaño de la muestra es lo suficientemente grande. Una transformación sugerida para remediar este problema es la de Freeman - Tukey que trata de hacer las desviaciones estandarizadas en cada celda cercanas a desviaciones estandar con media cero y varianza unitaria. Como es sabido

si se suman estas desviaciones al cuadrado se obtendrá una distribución χ^2 . La transformación de Freeman - Tukey se define:

$$z_i = \sqrt{x_i} + \sqrt{x_i + 1} - \sqrt{4\widehat{m}_i + 1}$$

que se utiliza para probar si el modelo es bueno así como para de tectar observaciones aberrantes. Cada una de las z, tiene - aproximadamente una distribución Normal estándar y como tal se puede utilizar para ver que celdas son aberrantes.

IV.D) Pruebas acerca de los parámetros.

Las estimaciones de M-V se obtienen al maximizar el kernel de la log-verosimilitud, éstas estimaciones se encuentran mediante el método iterativo que se ha presentado. La estimación $\widehat{\theta}$ está distribuída asintóticamente Normal y la correspondiente matriz de varianza-covarianza es la negativa de la inversa de la matriz Hessiana; esto es debido a que se trata de una estimación de M-V.

Dividiendo cada estimación por su correspondiente desviación estándar estimada se obtienen las estadísticas t correspondientes a cada parámetro.

Las hipótesis acerca de los parámetros en conjunto se establecen mediante pruebas de cociente de verosimilitud para muestras grandes. Este cociente de verosimilitud se define como:

$$\lambda = - 2 \log \frac{L(\widehat{\theta}_0)}{L(\widehat{\theta}_1)}$$

en donde $L(\widehat{\theta}_0)$ y $L(\widehat{\theta}_1)$ son las funciones de verosimilitud bajo los parámetros $\widehat{\theta}_0$ y $\widehat{\theta}_1$ respectivamente. En este caso $\widehat{\theta}_0$ involucra parámetros extra a los impues-

tos bajo $\widehat{\theta}_1$. λ se distribuye como una χ^2 bajo la hipótesis nula que θ es igual a $\widehat{\theta}_0$. El valor de λ es el valor de la prueba G^2 ($2 \mid 1$) con el modelo l (θ_1) anidado en el modelo 2 (θ_0), esto es para comparar dos modelos pero si se quiere probar la bondad de ajuste, se puede suponer que θ_0 es el vector que contiene a todos los parâmetros posibles, definiendo así un modelo satu-

X ;

Para comparar dos modelos, otra estadística que puede utilizarse es:

m:

rado, y entonces

$$R^{2} = \frac{G^{2}(1) - G^{2}(2)}{G^{2}(1)}$$

que es análoga a la R² utilizada en regresión y mide el éxito del modelo 2 al predecir la probabilidad.

IV.E) Una Medida de Asociación.

Goodman y Kruskal (1954) propusieron una medida de asociación entre dos variables categóricas ordenadas (por ejemplo, variables que presentan las categorías alto, normal, bajo, o bien crece, constante, decrece; etc.). Esta medida nos da, una idea de la magnitud de la asociación y de su dirección. La dirección guarda relación con la concentración de individuos en las celdas, por ejemplo, considérense las tablas siguientes para el caso 3 x 3:

× _{II}	0	0		0	0	X ₁₃
0	×22	0	у	0	× _{2 2}	0
0	0	×33		×3;	0	0

la primera tabla muestra asociación positiva perfecta mientras que la segunda muestra asociación negativa perfecta. La medida propuesta toma por definición los valores de +1 y -1 en cada caso respectivamente.

La base del coeficiente gamma de Goodman - Kruskal es un modelo probabilístico de comportamiento que se define a continuación.

Supóngase que dos individuos se seleccionan independientemente y en forma aleatoria de una población (con reemplazo, lo cual es innecesario en una población lo suficientemente grande). Cada uno de ellos cae en una celda de la tabla. Sea (i_1 , i_2) la celda en donde cae el l^2 individuo y (i_1 , i_2) la del segundo.

Si hay independencia se espera que el orden de las celdas no tenga conexión; si hay asociación positiva se espera que los individuos se concentren en la diagonal indicada, esto es, si $i_1 > i_1'$ entonces $i_2 > i_2'$, o si $i_1 < i_1'$ entonces $i_2 < i_2'$, a esto se le llama órdenes iguales. En caso de asociación negativa se espera que si $i_1 > i_1'$ entonces $i_2 < i_2'$, o si $i_1 < i_1'$ entonces $i_2 > i_2'$, lo que se conoce como órdenes diferentes. Hay empate cuando $i_1 = i_1'$ ó $i_2 = i_2'$ en este caso no se puede decir que tipo de relación hay.

Sean:

$$\Pi_{s} = \Pr \{ i_{1} > i'_{1} \ e \ i_{2} > i'_{2} \ \delta \ i_{1} < i'_{1} \ e \ i_{2} < i'_{2} \}$$

$$\Pi_{d} = \Pr \{ i_{1} > i'_{1} \ e \ i_{2} < i'_{2} \ \delta \ i_{1} < i'_{1} \ e \ i_{2} > i'_{2} \}$$

$$\Pi_{d} = \Pr \{ i_{1} = i'_{1} \ \delta \ i_{2} = i_{2} \}$$

entonces la probabilidad condicional de ordenes iguales dado que no hay empates es Π_s / (1 - Π_t) y la probabilidad condicional de ordenes diferentes dado que no hay empates es Π_d /(1 - Π_t).

Là gamma de Goodman - Kruskal se define como:

$$\gamma = \frac{\Pi_s - \Pi_d}{1 - \Pi_t}$$

que es la diferencia entre las probabilidades condicionales de ordenes iguales y diferentes. Mide que tanto más probable es obtener or denes iguales que diferentes en la tabla bivariada.

Como $\Pi_s+\Pi_d+\Pi_t=1$ entonces se pueden simplificar los cálculos para incluir sólo a Π_s y Π_t . Haciendo:

$$\gamma = \frac{\Pi_{s} - \Pi_{d}}{1 - \Pi_{t}} = \frac{\Pi_{s} - [1 - \Pi_{s} - \Pi_{t}]}{1 - \Pi_{t}} = \frac{2\Pi_{s} - 1 - \Pi_{t}}{1 - \Pi_{t}}$$

con:

$$\Pi_{s} = 2 \sum_{i_{1}} \sum_{i_{2}} p_{i_{1}i_{2}} \left\{ \sum_{i'_{1} > i_{1}} \sum_{i'_{2} > i'_{2}} p_{i'_{1}i'_{2}} \right\}$$

y

$$\Pi_{t} = \sum_{i_{1}} p_{i_{1}}^{2} + \sum_{i_{2}} p_{i_{1}} - \sum_{i_{1}} \sum_{i_{2}} p_{i_{1}}^{2}$$

Las propiedades de esta medida son las siguientes:

- i) γ es indeterminada si la población cae entera mente en un sólo renglón o columna de la tabla de clasificación cruzada.
- ii) γ es l si la población está concentrada en la diagonal arriba izquierda a abajo derecha de la tabla . γ es -l si la población está concentrada en la diagonal abajo izquierda a arriba derecha de la tabla.
- iii) γ es cero en el caso de independencia, pero lo contrario no necesariamente se cumple excepto para el caso 2 x 2 . Para poder establecer propiedades estadísticas de γ es necesaria una expresión para la varianza.

Esta medida sólo se aplica a tablas bivariadas de contingencia, cuando se estudian tablas multivariadas hay una generaliza ción debida a Kawasaki (1980) para medir asociación parcial, es decir, se mide la relación controlando la influencia de variables adicionales. A esta medida se le ha llamado la gama-componente γ_c .

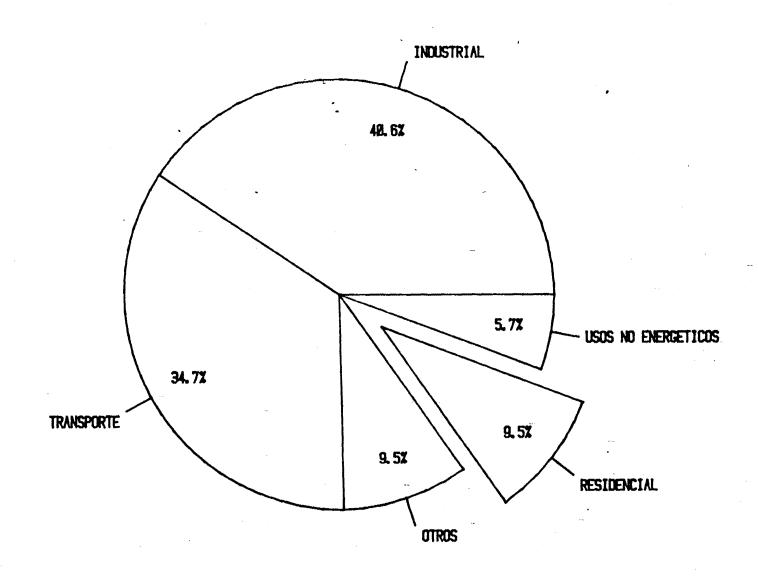
V. UNA APLICACION DEL MODELO

Introducción.

Se ha presentado el modelo log-lineal para analizar una tabla de contingencia. El objetivo de esta sección es plantear aquel modelo que resulte más adecuado para los datos bajo estudio, que presente las principales interrelaciones que existen entre las variables consideradas.

Dada la importancia que ha cobrado recientemente el estudio del consumo de energía, y como el sector residencial representa aproximadamente el 9% del gasto total en energía, (ver gráfica l), el objetivo del análisis es ver qué relación guardan entre si variables como el ingreso de la familia y su gasto en energía, si las familias con mayor ingreso gastan más energía, si los habitantes de zonas rurales se las arreglan con menos energéticos, si existe alguna relación entre el consumo de energía y el tamaño que tiene la vivienda, etc.

PARTICIPACION DE LA DEMANDA RESIDENCIAL EN LA DEMANDA TOTAL DE ENERGIA (1976)



A partir de estas interrogantes se pueden formular hipótesis de trabajo las cuales deberán probarse a partir de los datos. El objeto no es presentar al modelo log-lineal como la única posibilidad de análisis, sino señalar sus posibilidades y sus limitaciones. Evidentemente existen otras metodologías que, dependiendo del acceso a la información original, podrán arrojar mejores resultados o resultados complementarios; se pueden mencionar el análisis de varianza, modelos de regresión, detección automática de interacciones (AID), Analisis de Conglomerados (Cluster Analysis), etc.

Otro punto que conviene señalar es el relativo a la consideración del diseño muestral utilizado para obtener los datos. Exis te una controversia entre dos grandes escuelas de estadísticos.

Una señala que el diseño específico utilizado (estratificado, por conglomerados, etc.) debe incorporarse explícitamente en el análisis; la otra sostiene que esto no es necesario y que los análisis pueden llevarse a cabo como si se tratara de una muestra aleatoria. Para un comentario más a fondo, véase Tomberlin (1979 y 1980)

Imrey y Sobel (1980) y Fay (1982). Evidentemente este problema existe tanto en relación con la aplicación del modelo log-lineal a datos de encuesta, como con la aplicación de cualquier otro método de análisis: ANOVA, etc.. Si se tienen los datos originales y una descripción detallada del diseño muestral es posible intentar

formas de análisis más completas, sin embargo es bien sabido que con mucha frecuencia lo único que tiene es un cuadro y una idea vaga de cómo se obtuvieron los datos.

Así pues, el análisis que a continuación se presenta, supone una muestra aleatoria y como tal debe tomarse como un primer análisis de los datos.

V.1 Datos Utilizados.

Los datos utilizados corresponden a la Encuesta de Ingreso-Gasto de los Hogares (S.P.P. 1976). El paquete utilizado para dividir las variables en categorías fue el S.P.S.S. (Statistical Package for the Social Science, versión 6, 1º abril 1975), considerándose las siguientes variables con sus respectivas categorías:

1) Zona
$$\begin{cases} Rural & (i_1 = 1) \\ Urbana & (i_2 = 2) \end{cases}$$

La zona rural corresponde a municipios con menos de 100,000 habitantes y son las unidades primarias no autorepresentados (UPNAR) en la encuesta.

2) Número de cuartos
$$\begin{cases} 1-2 & (i_2=1) \\ 3-más & (i_2=2) \end{cases}$$

Para distinguir viviendas chicas y grandes. En esta variable, al igual que en las tres siguientes, se fijaron categorías hasta cierto punto arbitrarias, con el fin de aplicar los modelos log-lineales. Es bien conocido el hecho de que la elección de categorías influye en los resultados, Agresti (1976), por lo que nuestros resultados son condicionales en las categorías utilizadas.

3) Tamaño de la familia
$$\begin{cases} 0 - 3 \text{ hijos} & (i_3 = 1) \\ 4 - \text{más} & (i_3 = 2) \end{cases}$$

Estas categorías reflejan familias pequeñas y numerosas, ya que el número promedio de hijos es de aproximadamente 3.5 por pareja.

4) Ingreso (pesos/mes)
$$\begin{cases} 0 - 2400 & (i_4 = 1) \\ 2400 - 5725 & (i_4 = 2) \\ 5725 - mas & (i_4 = 3) \end{cases}$$

El ingreso así clasificado sirve para diferenciar ingresos bajos, medios y altos; estos niveles se consideran tomando en cuenta el salario mínimo vigente en esa época.

5) Gasto en energéticos (
$$i_5 = 1$$
)
(pesos/mes) $(i_5 = 1)$
($i_5 = 1$)
($i_5 = 2$)
($i_5 = 3$)

Esta categorización es un poco arbitraria pero intenta reflejar gastos bajos, medios y altos en energéticos. Para justificar esta clasificación puede observarse la tabla I que presenta para cada estrato de ingreso cual es su promedio de gasto en energéticos: para un nivel bajo de ingreso (0 - 2400) el gasto en energéticos fluctua entre 10.29 y 72.87, por lo tanto puede tomarse como referencia 50.00 que refleja de alguna manera un gasto medio aproximadamente. De esta manera pueden justificarse las otras divisiones.

PORCENTAJE DE GASTO EN ENERGETICOS
CON RESPECTO AL INGRESO

	****	. :
Estrato de Ingreso mensual	total de ene <u>r</u> géticos	gasto en ener- géticos*
0.1 - 700	2.94	10.29
700 - 1000	2.86	24.31
100 - 1350	2.80	32.9
1350 - 1800	3.22	50.71
1800 - 2400	3.47	72.87
2400 - 3150	3.89	107.95
3150 - 4300	3.66	136.34
4300 - 5725	3.89	194.99
5725 - 7500	4.35	287.64
7500 - 10150	4.57	403.30
10150 - 13400	5.32	626.43
13400 - 18000	4.64	728.48
18000 ó más	4.03	

Fuente: Encuesta Nacional de Ingresos y Gastos de los Hogares. 1977, S.P.P.

.

Tomando el valor medio del intervalo como referencia.

TABLA 2

									· ·
_	i ₁	1	l	l	1	2	2	2	2
	i ₂	- 1	l	2	2	1	1	2	2
	і з	1	2	1	2	l	2	1	2
14	i ₅					1	•		
1	1	602	1178	170	425	151	136	61	37
. 1	2	103	228	83	190	105	132	92	83
1	3	15	25	21	38	15	19	23	26
2	1	47	207	18	142	115	168	52	69
2	2	54·	251	67	360	170	507	211	601
2	3	11	80	36	194	30	177	95	395
3	ļ	9	19,	8	22	18	27	39	40
3	2	16	52	35	150	55	155	195	600
3	3	11	31	42	277	21	122	2 37	133

Variables consideradas

- l) Zona
- 2) Número de cuartos
- 3) Tamaño de la familia
- 4) Ingreso
- 5) Gasto en energéticos

Los datos utilizados en este capítulo se presentan en la Tabla 2, en donde aparecen los valores en las celdas elementales. La tabla consta de 72 celdas que resultan de considerar todas las posibles combinaciones de las categorías (2x2x2x3x3). El número total de familias en la muestra es de 11539. No existe el problema de celdas vacias debido a la magnitud de la muestra y gracias a esto fue posible estimar cualquier tipo de modelo.

En las secciones que siguen se presentan los resultados obtenidos al ajustar modelos log-lineales, mediante el programa CALM (Conditional Anova Log-Linear Models, desarrollado por John P. Link, Northwestern University. Noviembre, 1980), que permite ajustar modelos jerárquicos conjuntos y que presenta entre los resultados que imprime:

- a) Algunas medidas de bondad de ajuste: X^2 , G^2 , etc.
- b) Un análisis paramétrico: para cada configuración incluída en el modelo se presentan las estimaciones de los parámetros y su correspondiente t, así como la gama componente de Kawasaki para cada configuración bivariada.

Mediante estos resultados se analizan los modelos y el objetivo de esta sección será encontrar aquel modelo que presente un buen ajuste y que comparado con los demás sea el más adecuado.

V.2 Resultados obtenidos.

V.2A) Algunos modelos jerárquicos.

Para facilitar la comprensión se utilizarán las siguientes abreviaturas al referirse a las variables:

- 1) ZONA (Zona)
- 2) NCUAR (Número de cuartos)
- 3) TFAM (Tamaño de la familia
- 4) INGR (Ingreso)
- 5) GASTO (Gasto en energéticos)

Con objeto de examinar la relación que guardan entre sí estas variables, se consideraron los siguientes modelos:

1) Modelo saturado: que incluye a todos los efec-

tos posibles entre las variables, involucra a la configuración C_{12345} como estadística suficiente.

- Modelo con efectos de orden 4 : incluye to-dos los efectos hasta los de orden 4; las estadísticas suficientes son C_{1234} , C_{1235} , C_{1245} y C_{2345} .
- Modelo con efectos trivariados: cuyas configuraciones mínimas suficientes son: C_{123} , C_{124} , C_{125} , C_{134} , C_{135} , C_{145} , C_{234} C_{235} , C_{245} , donde los efectos incluídos involucran relaciones entre las variables indicadas en cada configuración.
- Modelo con efectos bivariados: que tiene como configuraciones mínimas C $_{12}$, C $_{13}$, C $_{14}$, C $_{15}$, C $_{23}$, C $_{24}$, C $_{25}$, C $_{34}$, C $_{35}$, C $_{45}$.
- 5) Modelo con efectos univariados: solamente incluye como configuraciones: C_1 , C_2 , C_3 , C_4 y C_5 . Este supone que no existe re-

lación alguna entre las variables.

En la Tabla 3 se presentan los valores de la X^2 , G^2 , los grados de libertad y el nivel de significancia asociado al ajuste de cada modelo (en parêntesis). Como la hipótesis nula a probar es lo de que el modelo no ajusta bien, una G^2 significativa, la cual corresponde a un valor menor a .1 entre parêntesis, implica que el modelo correspondiente no es bueno. Se observa que en los modelos H_4 y H_5 la G^2 es altamente significativa. Esto es una indicación de que no basta con incluír sólo efectos principales e interacciones de orden dos, síno que habrá que incluir al menos algunas interacciones de orden tres.

Formalmente esto se puede hacer mediante el uso de las G^2 , ya que éstas nos permiten decidir si los parámetros incluídos en un modelo y no en el otro hacen una contribución significativa al ajuste del modelo. Así, entre el modelo H_1 y el H_2 la diferencia está en que el primero incluye todas las interacciones, hasta las de orden cinco (saturado) mientras que el segundo sólo incluye hasta interacciones de orden cuatro; entonces G^2 (1/2)= G^2 (1) - G^2 (2) = 1.62 tiene una distribución con 4 grados de libertad. La comparación con el valor de tablas, χ^2_4 (.05) = 7.779 indica que la contribución no es significativa. La Tabla 4 presenta

TABLA 3

	suficientes	x ²	G ²	g.1.
H ₁	saturado	0.0 * (1.0)	0.0* (1.0)	0.
H ₂	de orden 4	1.63 (.804)	1.62 (.805)	4
H ₃	Trivariados	16.7 (.671)	16.8 (.668)	20
H ₄	Bivariados	216 (0.0)	201 (0.0)	45
H ₅	Univariados	18700 (0.0)	12500 (0,0)	66

^{*} Valor obtenido 0.0

las comparaciones entre los modelos H₂ y H₃ y la de H₃ con H₄. El valor entre paréntesis corresponde al nivel de significancia.

TABLA 4

	G 2	g.l.	
H ₃ - H ₄	184.2	2 5	
	(.005)		
н ₂ - н ₃	15.18	1 6	
	(.5115)		

Se tiene pues que la diferencia entre H_2 y H_3 no es significativa, es decir, las interacciones de orden 4 no son relevantes. Por el contrario, la diferencia de H_3 con H_4 es altamente significativa lo que implica que por lo menos algunas interacciones de orden tres sí son relevantes. Esto significa que se intentará encontrar un modelo entre H_3 y H_4 .

V.2B) Selección del Modelo.

A partir de un modelo que incluya los efectos trivariados se probará excluyendo aquellos parámetros menos significativos, hasta llegará aquel que posea un buen ajuste y un número menor de parámetros.

Teniendo como base el modelo con efectos trivariados se prueba excluyendo algunos de éstos parámetros que presentan menor significancia, de acuerdo a los valores χ^2 correspondientes. Las interacciones de orden 3 se denotan con la letra griega delta, δ .

En la Tabla 5 se presentan los resultados de excluir algunas configuraciones trivariadas. En ella se señalan varios modelos en los que se han omitido diversas configuraciones de orden tres. Estas configuraciones corresponden a los parámetros de las interacciones de orden tres. Así, en el modelo H_A se ha eliminado la interacción entre las variables 1, 2 y 5. La estadística G^2 permite analizar la bondad de ajuste de cada uno de ellos, al compararla con un valor de χ^2 obtenido de tablas, para los grados de libertad correspondientes, y que aparecen en la última columna de la tabla. Nuevamente una G^2 significativa

TABLA 5

Modelo	Configuraciones excluidas H ₃	G ²	g. l.
H 3		16.8 (.668)	20
H A	^C 125	17.5 (.733)	22
НВ	C ₁₂₅ C ₂₃₅	19.1 (.748)	24
^H C	C ₁₂₅ , C ₁₃₄ , C ₂₃₅	27.3 (.395)	26
H D	C ₁₂₅ , C ₂₃₄ C ₂₃₅	22.4 (.668)	26
H E	C_{125} , C_{134} , C_{235} , C_{234}	30.5 (.340)	28
H _F	$^{\text{C}}_{123}$, $^{\text{C}}_{125}$, $^{\text{C}}_{134}$, $^{\text{C}}_{235}$, $^{\text{C}}_{234}$	39.9	30

TABLA 6

Prueba	Parámetro Extra a Probar	G2	g.1.
H ₃ - Н _А	δ ₁₂₅	1.04 (.5945)	2
H_A - H_B	δ ₂₃₅	1.6 (.4492)	. 2
H _В - H _С	$\delta_{f 134}$	8.2 (.0166)	2
н _в - н _р	$^{\delta}$ 234	3.3 (.1921)	2
$H_D - H_E$	δ ₁₃₄	8.1 (.0174)	2
H _E - H _F	δ ₁₂₃	9.4 (.0091)	2

implica que el modelo correspondiente no es bueno. El nivel de significancia aparece entre parêntesis. Observamos que todos los modelos ajustan bien con excepción del $H_{\rm F}$, en el cual se ha eliminado el mayor número de interacciones. Como nos interesa obtener un modelo con el menor número de parámetros, parece ser que cualquiera en que se han eliminado 3 δ 4 interacciones podría ser adecuado.

Ahora veamos cuales de los parâmetros específicos hacen una contribución significativa al ajuste del modelo. Para esto pasamos a la Tabla 6. En ella utilizamos la propiedad de que G^2 (A/B) = G^2 (A) - G^2 (B) como se hizo al comparar modelos en la Tabla 3. Nótese que al referirse en la Tabla 6 que se está probando el parámetro δ_{125} en realidad se están probando los 12 parámetros correspondientes a las combinaciones de los niveles de las variables 1, 2 y 5: δ_{125} (i, j, k), i = 1,2 j = 1,2 y k = 1,2,3. Los dos grados de libertad se deben a que, con las restricciones tipo ANOVA, al fijar dos cualesquiera de los parámetros, los restantes (diez) quedan fijos. Entonces se tiene que:

- 1) H_A es preferible a H_3 , o sea que la eliminación de los parámetros δ_{125} no reduce en forma significativa el ajuste del modelo.
- 2) H_A es preferible a H_B , o sea que la interacción $\delta_{235} \quad \text{si contribuye significativamente al ajuste del modelo.}$
- 3) H_B es preferible a H_C ; no se puede excluir la configuración (interacciones de orden tres) C_{134} del modelo B sin afectar su ajuste.
- 4) H_D es preferible a H_B ; la combinación de interacciones C_{234} se puede eliminar del modelo B.
- Con base en lo anterior, se eliminan simultáneamente las interacciones que se indica en el modelo $H_{\bf E}$ y se compara con el $H_{\bf D}$ resultando que $H_{\bf D}$ es preferible a $H_{\bf E}$. No se puede eliminar la configuración C_{134} .
- 6) El modelo H_D es preferible al H_3 ya que G^2 (D/3) = G^2 (D) - G^2 (3) = 22.4 - 16.8 = 5.6 y tiene

6 g.1.

- 7) Se hace luego el intento de eliminar más interacciones de orden 3; por ejemplo la C₁₂₃, y esto resulta en el modelo F. Aquí la conclusión es que el modelo E es preferible al F.
- 8) El modelo finalmente seleccionado fue el H_D , del cual se excluyen las siguientes interacciones:
 - i) ZONA x NCUAL x GASTO
 - ii) NCUAR x TFAM x INGR
 - iii) NCUAR x TFAM x GASTO

No se consideró seguir eliminando parámetros ya que los demás efectos trivariados son significativos, y aunque el efecto bivariado β_{23} no es significativo, no es posible eliminarlo dada la estructura jerárquica implícita. Un intento de eliminarlo, es mediante el modelo H_F que excluía a la configuración C_{123} , pero este modelo se rechazó en favor de H_F .

V.2C) Interpretación de los resultados.

El modelo seleccionado en función de los parámetros que contiene es:

$$\ell_{\text{milizi3i4i5}} = \mu + \alpha_{1} + \alpha_{2} + \alpha_{3} + \alpha_{4} + \alpha_{5}$$

$$+ \beta_{12} + \beta_{13} + \beta_{14} + \beta_{15} + \beta_{23} + \beta_{24} + \beta_{25} + \beta_{34} + \beta_{35} + \beta_{45}$$

$$+ \delta_{123} + \delta_{124} + \delta_{i35} + \delta_{145} + \delta_{245} + \delta_{345} + \delta_{134}$$

El cual contiene 8 efectos menos que el modelo saturado. Puede utilizarse

$$G^2(1) - G^2(D) / G^2(1) = 12500 - 22.4 / 12500 = .9986$$

o sea que el 99.9% de la variación se explica mediante este modelo, lo cual representa un buen ajuste.

Se puede ver de la Tabla 7 que en general los valores observados y los estimados son muy semejantes, lo cual no es de sorprender dado el buen ajuste del modelo. Por otra parte los valores de la estadística Freeman-Tukey no muestran ningún comportamiento anormal; el valor más alto, 1.40, está dentro de los límites razonables para una normal estándar.

TABLA 7

		i ₁	1 1 1	1 1 2	1 2 1	i 2 2	2 1 1	2 i 2	2 2 i	
14	15	I3				4 .		<u> </u>		
1	1		602.0	1 178.0	170.0	425.0	151.0	Same	(1,0	;~,i
			591.06 (.45)	1 185.38 (21)	172.21 (15)	426.34 (05)	164,45 (-1,16)	12 5.10 (15 m)	59.25 (.05)	55 . 17 (* . 15)
1	2		103.0	224.0	53.0	190.0	105.0	132.6	92.0	\$2.0
			114.08 (-1.04)	220.63 (.51)	79.5 (.42)	169.79 (.03)	104.17 (.11)	(.27)	85,25 (,74)	ગરું .47 (સં.(ઝ)
1	3		15.0	25.0	21.0	35.0	15.0	19.0	23.0	25.0
•			16.13	23.71	21.02	38.14	14.56	19.55	22.9	26.55
			(22)	(.31)	(04)	(.02)	(.18)	(-,(%)	(.20)	(**(**)
2 :	1		47.0	207.0	18.0	142.0	115.0	168.0	52.0	13.1
			50.11	213.94	23.91	126.04	104.14	168.82	53.84	77.04
			(41)	(46)	(-1.23)	(1.40)	(1.06)	((∺)	(22)	(۶١)
2	2		54.0	251.0	67.0	360.0	170.0	9,7	201.0	6
			52.08	24 8.9 (.15)	62,47 (.59)	368.55 (43)	168.51 (.13)	.512.51 (23)	2.5,24 15,52)	554. 4 (.50)
			(.30)	(+12)		•				
2	3	•	11.0	80.0	36.0	194.0	30.0	77	15.0	35
			13.38 (60)	71.59 (.99)	31.(4 (.89)	204.98 (76)	36.32 (-1.03)	171 .71 (559)	4, .25 (.41)	352.72 ()
			(00)	(.77)	(.03)	(*.70)	(-1.ω)	(+37)	1.74)	:
3	1		9.0	19.0	8.0	22.0	15.0	27.11	70.0	40.
			8.48	18.75	8.23	22.51	20.26	25,49	5.14	41.23
			(.25)	(.11)	(0.0)	(06)	(46)	(404) 	(.55)	(15)
3	2		16.0	52.0	35.0	150.0	55.0	255.0	195.0	. 3 6.76
•	•		15.39	53.93	34.48	149.2	48.11	20.57	2.52	5 9
			(.21)	(2 3)	(.13)	(.09)	(.99)	(42	(55)	(.25)
3	3		11.0	31.0	42.0	2 77.0	21.0	122.0	257.0	1 307.6
			7.29	34.13	47.13	272,44	19.45	124.10	237.10	1 335. 2
			(1.29)	(5 0)	(~.73)	(.29)	(.39)	(- ,17)	· (.(.)	(.(0)

Valor Observado Valor Ajustado (componente Freeman-Tukey)

TABL'A 7

		i ₁ i ₂ i ₃	1 1	1	1 2	1 2	2 1	2 1	2 2	2 2
i 4	i ₅	i ₃	i	2	1	2	ī	2	ī	2
1	1		602.0	1 178.0	170.0	425.0	151.0	136.0	61.0	37. 0
-	-		591.06	1 185.38	172.21	426.34	164.45	126.10	56.28	38.17
			(.45)	(21)	(15)	(05)	(-1 .05)	(.89)	(.65)	(-,15)
1	2		103.0	228.0	83.0	190.0	105.0	132.0	92.0	83.0
			114.08	220.63	79.5	189.79	104.17	129.11	85.2 5	93.47
			(-1.04)	(.51)	(.42)	(.03)	(.11)	(.27)	(.74)	(4.09)
1	3		15.0	2 5.0	21.0	38.0	15.0	19.0	23.0	26.0
			16.13	23.71	21.02	38.14	14.56	19.66	22.9	26.55
			(22)	(.31)	(04)	(.02)	(.18)	(08)	(.20)	(06)
2	1		47.0	207.0	18.0	142.0	115.0	168.0	52.0	69.0
			50.11	213.94	2 3. 9 1	126.04	104.14	168.82	53.84	7 7.2
			(41)	(46)	(-1.2 3)	(1.40)	(1.06)	(-,04)	(22)	(93)
2	2		54.0	2 51.0	67.0	360.0	170.0	507.0	211.0	601.0
			52.08	248.9	62.47	368.55	168.51	512.51	218.94	589.04
		,	(.30)	(.15)	(.59)	(43)	(.13)	(2 3)	(52)	(.50)
2	3		11.0	80.0	36.0	194.0	30.0	177.0	95.0	395.0
			13.38	71.59	31.04	204.98	36.32	176.71	91.25	392.72
٠			(60)	(.99)	(.89)	(76)	(-1.05)	(.59)	(.41)	(.13)
3	1		9.0	19.0	8.0	22.0	18.0	27.0	39.0	40.0
			8.48	18.78	8.23	2 2.51	20.26	25.49	37.04	41.22
			(.25)	(.11)	(0.0)	(~.06)	(- .4 6)	(.34)	(.36)	(15)
3	2		16.0	£9 0	35.0	150.0	55.0	155.0	10F 0	404.0
3	2		15.39	52. 0 53.9 3	34.48	150.0 149.2	55.0	155.0	195.0	606.0
			(.21)	(23)	(.13)	(.09)	48.11 (.99)	160.57 (42	203.02	599.3
			(-41)	(~.43)	(.13)	(+07)	(.77)	(=2	(5 5)	(.28)
3	3		11.0	31.0	42.0	277.0	21.0	122.0		1 337.0
			7.29	34.13	47.13	272.44	19.48	124.10	237.10	1 336.32
			(1.29)	(50)	(- . 73)	(.29)	(.39)	(17)	(.01)	(.03)

Valor Observado Valor Ajustado (componente Freeman-Tukey) La Tabla 8 presenta los valores de los parâmetros estimados para el modelo, correspondientes a los efectos principales y a las interacciones de orden dos. Por la dificultad en interpretar las de orden mayor éstas sólo se presentan en forma de listado en un apéndice al final del artículo. Junto con los parâmetros correspondientes a las interacciones se ha incluído la gama - componente que mide el grado de asociación entre las dos variables a que se refiere la interacción, eliminando el efecto de los demás. Conviene señalar que la gama - componente sólo tiene sentido en el caso de que ambas variables involucradas sean ordinales, cosa que es válida en todas las variables, pues aún la distinción entre zona rural y urbana implica un orden en cuanto al tamaño de la localidad.

A partir de esta tabla se observa que todos los efectos principales son significativos; esto se determina a partir de los valores t que aparecen entre paréntesis. Esto quiere decir que cada una de ellas contribuye a determinar el número de personas que cae en cada celda. Los resultados para las interacciones indican que con excepción de ZONA x NCUAR y NCUAR x TFAM todos son significativos. La gama-componente de los restantes son significativos y positivos lo cual es de esperarse. La interacción ZONA x TFAM es negativa, indicando que hay más familias pequeñas en la zona urbana y más familias numerosas en la zona rural.

			ZONA x NCUAR:				NCUAR x INGR:			
Rural 1849 (10.03)	Urbana .1849 (10.63)		Rural	Chica (1-2) .0311 (2.161)	Grande (>2)(3]] (2.16])		Chica	Bajo ,3430 (14.621)	Medio .0943 (4.946)	Alta 3 (17.5%)
DE CUARTOS (NCUA	·		Urbana	0311 (2.161)	+ .0311 (2 .161)		Grande	3430 (14.621)	0943 (4.846)	.4373 (17.351)
Chica (de 1 a 2) U 2 i	Grande (3 y mús) 1725			GAMA COMP. =	.0622 (2.167)			GAMA COMP. =	.4777 (20.799)	
(15,153)	(10.133)		ZONA x TFAM:				NCUAR x GASTO:			
DE L'AMILIA (TEAM Pequeña (de 0 a 3 bijos)	Numerosa (4 o más bilos)		Rural	Pequeña (0-3) 1338 (9.508)	Numcrosa(>3) .1338 (9.508)		Chica	Bajo .4225 (18.346)	Media 0158 (.872)	Alto 7.43(7 (17.35%)
4759 (25.440)	.4759 (28.446)		Urbana	.1338 (9.508)	1338 (9.508)		Grande	4225 (18.346)	.0158 (.872)	.40/7 (11.157)
INGRESO (INGR): - 100	<u>Medio</u> -3531	Alto 3840		GAMA COMP. =	-, 2 613 (9,968)			GAMA COMP. =	.5632 (23.513)	
(1.250)	(15.739)	(13.820)	ZONA x INGRESO:				TFAM x INGRESO:			
GASTO (GASTO):	<u>Medio</u>	Alto	Rural	Bajo .5353 (6.975)	Medio 1512 (6.975)	Alto 4051 (14.680)	Pequefia	Bajo .2993 (12.499)	Madio 1:09 (6.945)	Alto - 12-4 (4,125)
7.12i0 (4.64s)	,5155 (24 ,171)	3944 (14.669)	Urbana	5563 (23.024)	.1512 (6.975)	.4051 (14.680)	Numerosa	2993 (12 .499)	.1473 (8.945)	.19:4 (1.ಟನ್)
				GAMA COMP. =	,5646 (26,209)		•	GAMA COMP. =	- ,2637 - (9,390)	
	e e e e e e e e e e e e e e e e e e e		ZONA x GASTO:		•		TFAM x GASTO:			
			Rural	Bajo .3078 (12.219)	Medio 1380 (6.834)	Alto 1698 (6.614)	Pequeña	Bajo .1824 (7.660)	Mllo -10359 (1.861)	(5.575)
			Urban a	3078 (12.219)	.1380 (6.834)	.1698 (6.614)	Numerosa	-,1824 (7,660)	+.0359 (1.861)	.1465 (5.875)
				GAMA COMP. ≈	.3072 (10.915)			GAMA COMP.	2!58 (7.577)	
			NCUAR x TFAM:				ingr x gasto:		-	, AN
			Chica (1-2)	<u>Pequeña</u> (0-3) .0110 (.888)	Numerosa (>3) 0110 (.888)		Bajo	Bajo 1.0226 (31.465)	Medio 1179 (4.033)	<u>Alto</u> 5047 (22.213
			Grande (> 2)	0110 (\$28.)	.0110 (888.)		Medio	1598 (5.032)	.1103 (4.22°)	.0490 (1,451)
				GAMA COMP. ≈	.0220 (.888)		Alto	8ó32 (20.150)	.0075 (.232)	,515°° (23.1°°

Se pueden resumir los hallazgos a partir de las interacciones y los coeficientes gama - componente como sigue:

- i) el gasto en energía está áltamente relacionado con el nivel de ingreso.
- ii) el ingreso varía de acuerdo a la zona.
- iii) el gasto en energéticos y el número de cuartos están relacionados entre sí.
- iv) el número de cuartos que posee la vivienda y el nivel de ingreso están relacionados.
- v) el tamaño de la familia guarda poca relación con el ingreso y con el gasto, pues las gamas son pequeñas aunque significativas.
- vi) el gasto depende de la zona.
- vii) las familias de las zonas rurales tienen más hijos, mientras que las de las zonas urbanas menos.
- viii) el número de cuartos y el tamaño de la familia no tienen ninguna relación.

Otro tipo de análisis que se puede hacer a partir del modelo es con base en los "momios" para variables con dos categorías. Estos son el cociente de la probabilidad de caer en una categoría entre la de caer en la otra. Con base en el modelo se representa de la siguiente manera para la variable ZONA:

$$\frac{m_{1} i_{2} i_{3} i_{4} i_{5}}{m_{2} i_{2} i_{3} i_{4} i_{5}} = \frac{e^{\alpha_{1}(1)}}{e^{\alpha_{2}(2)}} \cdot \frac{e^{\beta_{12}(1, i_{2})}}{e^{\beta_{12}(2, i_{2})}}$$

$$\frac{e^{\beta_{13}(1, i_{3})}}{e^{\beta_{13}(2, i_{3})}} \cdot \frac{e^{\beta_{14}(1, i_{4})}}{e^{\beta_{14}(2, i_{4})}} \cdot \frac{e^{\beta_{15}(1, i_{5})}}{e^{\beta_{15}(2, i_{5})}} \cdot \frac{e^{\delta_{123}(1, i_{2}, i_{3})}}{e^{\delta_{123}(2, i_{2}, i_{3})}}$$

$$\frac{e^{\delta_{124}(1, i_{2}, i_{4})}}{e^{\delta_{124}(2, i_{2}, i_{4})}} \cdot \frac{e^{\delta_{135}(1, i_{3}, i_{5})}}{e^{\delta_{135}(2, i_{3}, i_{5})}} \cdot \frac{e^{\delta_{145}(1, i_{4}, i_{5})}}{e^{\delta_{145}(2, i_{4}, i_{5})}}$$

$$\frac{e^{\delta_{134}(1, i_{3}, i_{4})}}{e^{\delta_{134}(2, i_{3}, i_{4})}}$$

que se interpreta como el "chance" relativo de que una familia caiga en la categoría rural en lugar de en la urbana, dependiendo de los niveles de las otras variables.

El valor de los momios correspondientes se presentan en la Tabla 9 para cada uno de los valores de los índices para las variables TFAM, NCUAR, INGR, GASTO. Esta tabla se interpreta de la siguiente manera, por ejemplo si se observa en la Tabla el valor de la celda para $i_2=1$ $i_3=1$, $i_4=1$, $i_5=1$ se ve que hay aproximadamente 4 familias de la zona rural por una de la zona urbana que caen en esta celda (ya que el valor del momio es de 3.75).

De acuerdo a esta tabla si la familia presenta un ingreso medio superior es más probable pertenecer a la zona urbana que a la rural y tener un consumo de energía medio o superior. También se observa que es más probable que para una familia de la zona rural su patrón de gasto de energía sea menor (los momios más altos se observan en gasto bajo en energéticos y nivel de ingreso bajo). Otra característica que hacen resaltar los momios es que las familias de bajos ingresos se concentran en las zonas rurales, así como que en general el tamaño de la familia sea mayor en las zonas rurales (que en la tabla se observa ya que los momios de estar en la categoría 4 hijos ó más son en general mayores que los de estar en la categoría de 1 a 3 hijos para un mismo tamaño de vivienda, ingreso y gasto en energía).

En la Tabla 10 se presenta un análisis similar utilizando los momios de NCUAR. Se observa que el tamaño de la vivien da aumenta con el nivel de ingreso.

TABLA 9
MOMIOS PARA RURAL/URBANA

		NCUAR				
Gasto en	Ingreso	i,	= 1	i 2	= 2	
Energéticos		T F A M				
		i ₃ = 1	i ₃ = 2	i ₃ = 1	i ₃ = 2	
i ₅ = 1	i 4 = 1 i 4 = 2 i 4 = 3	3.59 0.481 0.467	9.407 1.268 0.690	3.059 0.444 0.222	11.182 1.635 0.547	
i ₅ = 2	$i_{4} = 1$ $i_{4} = 2$ $i_{4} = 3$	1.094 0.528 0.546	1.709 0.486 0.336	0.884 0.271 0.170	2.032 0.626 0.249	
i ₅ = 3	$ \begin{array}{ccc} $	1.108 0.369 0.374	1.211 0.406 0.275	0.943 0.340 0.199	1.439 0.734 0.204	

TABLA 10

MOMIOS DE TAMAÑO DE LA VIVIENDA CHICA /GRANDE

		ZONA				
Gasto en Energéticos	Ingreso	iı	= 1	i	1 = 2	
		T F A M				
		i ₃ = 1	i ₃ = 2	i ₃ = 1	i3 = 2	
i ₅ = 1	i 4 = 1 i 4 = 2 i 4 = 3	3.433 2.183 1.072	2.781 1.697 0.833	2.923 1.934 0.547	3.305 2.186 0.618	
i ₅ = 2	$i_{4} = 1$ $i_{4} = 2$ $i_{4} = 3$	1.435 0.869 0.446	1.162 0.676 0.361	1.222 0.770 0.237	1.382 0.870 0.268	
i ₅ = 3	$i_{4} = 1$ $i_{4} = 2$ $i_{4} = 3$	0.766 0.431 0.154	0.621 0.349 0.125	0.653 0.398 0.082	0.738 0.450 0.093	

TABLA 11

MOMIOS DE TAMAÑO DE LA FAMILIA PEQUEÑA/NUMEROSA

		ZONA				
Gasto en	Ingreso	$i_1 = 1$		$i_1 = 2$		
Energéticos		NCUAR				
	Ī	i ₂ = 1	i ₂ = 2	i ₂ = 1	i ₂ = 2	
_	$i_{4} = 1$ $i_{4} = 2$ $i_{4} = 3$	0.498	0.404	1.306	1.475	
$i_5 = 1$	i 4 = 2	0.236	0.191	0.617	0.698	
	$i_4 = 3$	0.451	0.365	0.794	0.898	
	i 4 = 1	0.517	0.418	0.807	0.912	
$i_{5} = 2$	$i\frac{4}{4}=2$	0.209	0.170	0.329	0.372	
5	$ \begin{array}{r} i_4 = 1 \\ i_4 = 2 \\ i_4 = 3 \end{array} $	0.285	0.231	0.300	0.339	
3		0.680	0.551	0.743	0.840	
i ₅ = 3	$\frac{1}{1}\frac{4}{4} = \frac{1}{2}$	0.187	0.152	0.206	0.233	
-5 ~	$i_{4} = 1$ $i_{4} = 2$ $i_{4} = 3$	0.214	0.173	0.157	0.178	

Al comparar los momios entre las zonas rural y urbana estos aumentan para las familias pequeñas o sea que en la zona urbana hay una tendencia al hacinamiento.

La Tabla II contiene los momios para el tamaño de la familia. En ella se nota que en general predominan las familias numerosas, con excepción de los niveles de ingreso bajo en la zona urbana. Los que más cerca le siguen son los niveles de ingreso alto en la misma zona.

En la zona urbana el número de familias numerosas aumenta con el ingreso, mientras que en la rural el aumento se da al pasar de niveles de ingreso bajos al nivel medio, pero disminuye nuevamente al pasar a ingresos altos. El número de cuartos no parece tener ningún efecto. Lo cual es consistente con lo señalado a partir de la Tabla 10.

V.2D, Conclusiones

Hemos presentado una aplicación del modelo log-lineal para el análisis de los datos sobre ingresos y gastos familiares, con especial interés en el gasto en energéticos. Hemos confirmado

algunas hipótesis como son el que el gasto depende del nivel de ingreso, del número de cuartos y de la zona, mas no del número de miembros en la familia. Esto es congruente con el hecho de que a ciertos niveles de ingreso existe hacinamiento en viviendas de uno y dos cuartos.

Algunas de las variables utilizadas son de tipo "intervalo", y se categorizaron para este estudio. Un análisis más profundo, y ya con la información que hemos obtenido con el presen
te podría llevarse a cabo modificando los criterios de clasificación o bien, si estuviera disponible la información original, utilizando alguna de las técnicas alternativas.

Hemos querido ilustrar la aplicación del modelo log-lineal a datos económicos y contribuir a su difusión como un método más dentro del análisis econométrico, concientes de que tiene
limitaciones, pero que en el caso de variables que por su naturaleza son categóricas tal vez sea lo más adecuado.

TABLA 8 (Continuación)

Parámetros Estimados del Modelo H_D Efectos Trivariados

•				_	R.			0	
		ı		ZONA	NCUAR	TFAM	8	GAST(
ı				Õ	Ď	Έ	8	¥	
								0	
				-		 -			
- 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4 1 4	. 22795-1	3.2 <i>r=</i>) -3.3 = 9 -3.3 = 9		<u>-</u>	<u>;</u>	1-7	<u>)</u> ()	· <u></u>	
	· · · · · · · · · · · · · · · · · · ·			넄. —	۲	$\frac{2}{1}$		<u> </u>	
		3.3.69 -3.3.69 -3.3.669		<u> </u>	<u> </u>		0 0 0	<u>C</u> _	
				-1	1	<u>-</u>	13	` <u>-</u>	
.61637				77	1	?	-5-		
		3.3169		· ·	2	-			
- 4151 ys - 1		-3.3669					<u>.</u>	_{(*}	
and a second sec					- ·				
programme and the second secon	The state of the s		1						
									-
	731-1	-1.754		-	-	1)		C	
• • • • • • • • • • • • • • • • • • • •	111-11	-3.116/		1	, ;	4.3			
• 1 2 2 2 2	. 1236E-)1	4.524		1		1.	2	C	
• 12 2 2 3 3 5 - 1	.145737-01	2.7:45	i.,	1	7	n]	ſ	
This is the first term of the first ter	·145775-2	1.1666					?	(:	
	,010365	- 3 6		1		1	? ? ? ? ? 1		
	- F ラヤキュー 1	1 . 7 . 44			<u>.</u>		1		
	110000	3.1-66) 	1	· .	2	. : _	
	· 12 / E- 1	-4.024		3	· ·		-3	. (,	
	. 11737-01	-1.7546		<u>:</u>	· \ 	. 1 ⁵	•		
and the second s	1. 1. 1. 7. 7. 7. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.	-3.1 +		<u>.</u>			2	<u>. </u>	
	- 1	1.6534		2	3	. • 	3		
	7	_:				-,	 -		
- 10 (4) (5) (5) (6) (7) (7) (7) (7) (7) (7) (7		-1.1.17 -1.64-1				-	<u> </u>	-,*	
				1		1	-		
						<u>.</u>	1	-,	
						4-	-		
		<u> </u>		1	 .	5	2	i	
	جو ايان مدي ميٽسيدي، يوان استران ايان ايان ايان	1, 175		- - -		•	-	<u>.</u>	
						1	- ;	ſ.	
AND THE REAL PROPERTY OF THE P		T-7,:54%			0	1	· ·	Ö	
		-1,7337		7	1.	2	1	<u>,</u>	
		F144 1 1 1 1 -1 -1 -1		7	1	<u> </u>	· · · · ·		
****					<u> </u>	÷			
	in program i							-	

TABLA 8 (Continuación)
Parámetros Estimados del Modelo H_D
Efectos Trivariados

<u></u>		i	ZONA NCUAR TFAM INGR GASTO
• 3 2 2		-5.47.2 .77.(75 4.3213	1 0 1 0 2
		• // (7)	
			1 0 1 0 3
	•		1
7		-4.3214	
		7.67-	1 0 2 0 3 2 0 3 C 2
		* 4 , 7	, 1 2
		-5.4771 -6.65.72	
			2 0 0
	المناسطين أدارات المائات	a managaran da ang ang ang ang ang ang ang ang ang an	
	VALUE OF STREET OF STREET	and the state of t	Contraction of the Contraction o
And the state of t		. 4 (4)	1
		-3.1633	1 3 7 1 3
7:7:17:18-		-0.45	3 (2)
	7.1	· Q 5. 1 × 3	2 2
		1, 7741	1 2 3
		3 . 1 2 5 5	1 0 0 3 1
	. • 14 / NATH-112 11	112,202	<u> </u>
24 2 22	المستعلق المراجع المرا		1 1 3 3
	. •		
and the second s	. 1117-1	7257	2 2 1 2 1 2
		<u> </u>	- <u> </u>
		25 3 4	2 2 2 1 -
			2
	-		2 1 3 2
	•		<u> </u>
processing and the state of the			

TABLA 8 (Continuación)
Parámetros Estimados del Modelo H_D
Efectos Trivariados

:			ZONA NCUAR TFAM INGR GASTO
. f			ZONA NCUAI TFAM INGR GASTO
-,274:27-1	9504E-01	93102	
251435-10	.25610E-01	945(7	0 1 0 1 2
.526177-11	.381516-01	1.3792	0 1 6 1 3
- 10250 5-4	.28173E-W	19173	(1 0 2 1
27-14'- 1	•:2169E- 1	-1.7573	$\frac{1}{\sqrt{2}}$
:332702mt,	.[8115E-01	1.1832	1 0 2 2
.325/43-03	1.385495-01	.85253	7 1 (3 1
. 3 3 3 3 2 - 1 3	.277405-01	1.9114	0 1 0 3 2
	20232-01	-2.682)] 1 0 3 3
. 2746.25-	.29:04:-01	.93102	0 2 0 1 1
.251408-11	.766)05-01	.94=07	0 2 6 1 2
		-1.0792	0 2 0 1 3
	• 11 1721 - 1	.19153	0 2 0 2 1
.27:745-12	.001/05-11	1.2573	0 2 6 2 2
— , ₹ ↑ 2 ↑ · · · · · · · · · · · · · · · · ·	• 219 (N.A.—.)	-1,1:33	<u> </u>
325 400	. 1431-01 . 17498-01	<u>05153</u>	0 3 1
 > 20200 a •	. Tare	-].C' ₁ 4	2 9 3 2
.258171-11	•086835 - -01	2.6820	0 2 0 3 3
12441	.,:11555	<u> </u>	0 0 1 1 1
17: 4581	. / [41][-0]	56604	0 0 1 1 2
.1414.	. (2575-0] . (2575-0] . (35, 25-0]	3.27 7	0 0 1 1 3
• 2258. t = 1		24511	0 0 1 2 1
- <u>• 4466 </u>	. 3,. 2,11	.94635	
123531-	. 14.5741	41247	
	<u> </u>	3.4607	
	.11 .71-01 .11191-01	-4.0037 3.9920	
. 17441		3.777	
			2 1 2
-, -, -, -	to the filter	.50004 -3.5007 .34511	<u>2</u>
	. 124621-01	0145.1	0 0 2 2 1
	** 236 × 1 = 01	94/29	0 (? / 2
			2 2 3
		-1,1417	<u> </u>
		.1978-	0 2 3 2
1750		-9.2407 -1976- 4.6597	
	and the second s		

BIBLIOGRAFIA

- Agresti, A., The Effect of Category Choice on Some Ordinal Measures of Association, Journal of the American Statistical Association, Vol. 71, pp 49 53.
- Bishop, Y.M.M., Fienberg, E.S., Holland, P.W. <u>Discrete</u>

 <u>Multivariate Analysis: Theory and Practice</u>. Ed. MIT Press.

 <u>Sexta Edición</u>, Cambridge, Massachusetts (1980).
- 3) Fay, R., Contingency Table Analysis for Complex Sample Designs: CPLX, Proceedings of the Survey Research Section ASA. (1982).
- 4) Haberman, Shelby. Analysis of Qualitative Data. Vol. I Introductory Topics. Ed. Academic Press. Nueva York. (1978).
- Imrey, P.B., y E. Sobel, Modeling Contingency Tables From Complex Surveys, Proceedings of the Survey Research Section ASA (1978).
- 6) Kawasaki, Seiichi. Application of Log-Linear Probability Models in Econometrics. Tesis doctoral presentada en la Universidad de Northwestern. Departamento de Economía (1979).
- 7) Link, John P. CALM. <u>Conditional ANOVA Log-Linear Models:</u>
 User's Guide. Universidad de Northwestern (1980).
- 8) Maxwell, A.E. Análisis Estadístico de Datos Cualitativos. Ed. UTEHA (Manuales 222) traducción al español por Ramón Galhd Garza. México (1965).
- 9) Mendoza, Y., El Modelo Log-Lineal y su Aplicación. Tesis Profesional para obtener el Título de Licenciado en Matemáticas Aplicadas (1984). (ITAM)
- Nerlove, M., Press, S.J. Univariate and Multivariate Log-Linear and Logistic Models. Preparado para la Administración del Desarrollo Económico y el Instituto Nacional de Salud. RAND CORPORATION R 1306 EDA&NIH. (diciembre, 1973).

- Nerlove, M., Press S.J. Multivariate Log-Linear Probability Models for the Analysis of Qualitative Data. Artículo para discusión Num. 1.

 Universidad de Northwestern. Centro de Estadística y Probabilidad (1976).
- Nerlove, Marc. <u>Expectations</u>, <u>Plans and Realizations in Theory and Practice</u>. Econométrica, Vol. 51. Num 5, (septiembre, 1983).
- Payne, C. "The Log-Linear Model for Contingency Tables". Tomado de The Analysis of Survey Data: Exploring Data Structures. Ed. John Wiley and Sons. Nueva York (1977) pags. 106.
- Samaniego, R., Berndt, E. Residential Energy Demand in México. publicado por el Instituto de Tecnología de Massachussetts. Laboratorio de Energía. (Agosto, 1982).
- Tomberlin, T.J., The Analysis of Contingency Tables of Data From Complex Surveys, Proceedings of the Survey Research Section, ASA (1979).
- Tomberlin, T.J., A Model Based Approach to the Analysis of Contingency Tables of Data From Complex Surveys, Proceedings of the Survey Research Section, ASA (1980).

Serie Documentos de Trabajo 1984

- No. I Alberro, José Luis, "Introduction and Benefit of Technological Change under Oligopoly".
- No. II Serra-Puche, Jaime y Ortíz, Guillermo, "A Note on the Burden of the Mexican Foreign Debt".
- No. III Bhaduri, Amit, "The Indebted Growth Process".
- No. IV Easterly, William, "Devaluation in a Dollarized Economy".
- No. V Unger, Kurt, "Las Empresas Extranjeras en el Comercio Exterior de Manufacturas Modernas en México".
- No. VI De Alba, Enrique y Yolanda Mendoza, "El Uso de Modelos Log-Lineales para el Análisis del Consumo Residencial de Energía".

Serie Documentos de Trabajo 1983

- No. I Bhaduri, Amit "Multimarket Classification of Unemployment".
- No. II Ize, Alain y Salas, Javier "Prices and Output in the Mexican Economy: Empirical Testing of Alternative Hypotheses".
- No. III Alberro, José Luis "Inventory Valuation, Realization Problems and Aggregate Demand".
- No. IV Sachs, Jeffrey "Theoretical Issues in International Borrowing"
- No. V Ize, Alain y Ortiz, Guillermo "Political Risk, Asset Substitution and Exchange Rate Dynamics".
- No. VI Lustig, Nora "Políticas de Consumo Alimentario: Una Comparación de los Efectos en Equilibrio Parcial y Equilibrio General".
- No. VII Seade, Jesús "Shifting Oligopolistic Equilibria: Profit-Raising Cost Increases and the Effects of Excise Tax".
- No. VIII Jarque, Carlos M. "A Clustering Procedure for the Estimation of Econometric Models with Systematic Parameter Variation".
- No. IX Nadal, Alejandro "La Construcción del Concepto de Mercancía en la Teoría Económica".
- No. X Cárdenas, Enrique "Some Issues on Mexico's Nineteenth Century Depression".
- No. XI Nadal, Alejandro "Dinero y Valor de Uso: La Noción de Riqueza en la Génesis de la Economía Política".
- No. XII Blanco, Herminio y Garber, Peter M. "Recurrent Devaluation and Speculative Attacks on the Mexican Peso".

El Centro de Estudios Económicos de El Colegio de México, ha creado la serie "Documentos de Trabajo" para difundir investigaciones que contribuyen a la discusión de importantes problemas teóricos y empíricos aunque estén en versión preliminar. Con esta publicación se pretende estimular el análisis de las ideas aquí expuestas y la comunicación con sus autores. El contenido de los trabajos es responsabilidad exclusiva de los autores.

Editor: José Luis Alberro

Serie Documentos de Trabajo 1982

No. I	Ize, Alair	"Disequilibrium Theories, Imperfect Compet-
	·	ition and Income Distribution: A Fix Price Analysis"

- No. II Levy, Santiago "Un Modelo de Simulación de Precios para la Economía Mexicana"
- No. III Persky, Joseph and Tam, Mo-Yin S. "On the Theory of Optimal Convergence"
- No. IV Kehoe, Timothy J., Serra-Puche, Jaime y Solis, Leopoldo
 "A General Equilibrium Model of Domestic
 Commerce in Mexico"
- No. V Guerrero, Víctor M. "Medición de los Efectos Inflacionarios Causados por Algunas Decisiones Gubernamentales: Teoría y Aplicaciones del Análisis de Intervención"
- No. VI Gibson, Bill, Lustig, Nora and Taylor, Lance "Terms of Trade and Class Conflict in a Computable General Equilibrium Model for Mexico"
- No. VII Dávila, Enrique "The Price System in Cantillon's Feudal Mercantile Model"
- No. VIII Ize, Alain "A Dynamic Model of Financial Intermediation in a Semi-Industrialized Economy"
- No. IX Seade, Jesús "On Utilitarianism and Horizontal Equity:
 When is the Equality of Incomes as such
 Desirable?"
- No. X Cárdenas, Enrique "La Industrialización en México Durante la Gran Recesión: Política Pública y Respuesta Privada"