

Borrador
Enero, 1977

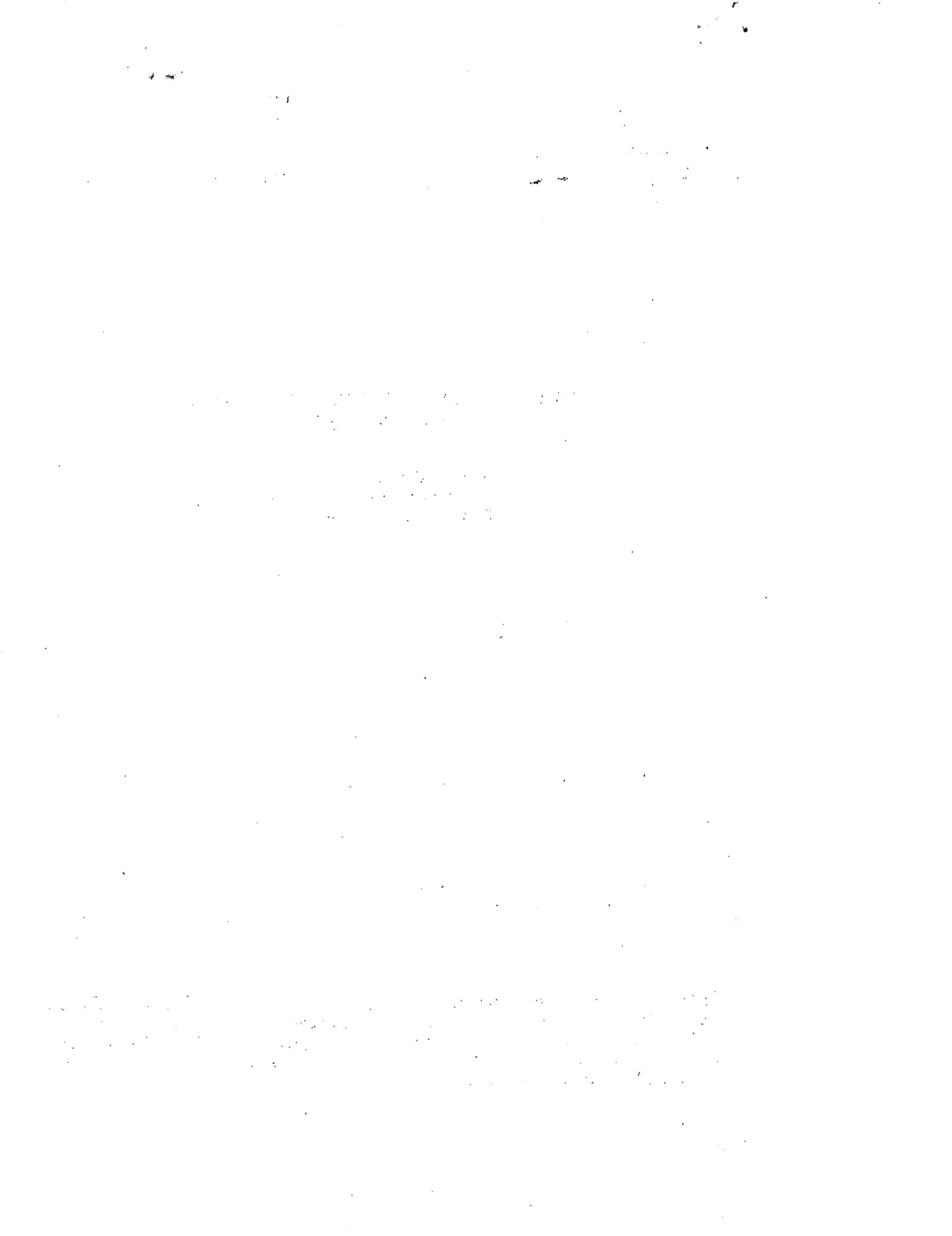
Comisión Económica para
América Latina

Banco Mundial
Development Research Center

ANALISIS DE DESCOMPOSICION: UNA GENERALIZACION
DEL METODO DE THEIL II/

Oscar Altimir
Ana Crivelli
Sebastián Piñera

II/ Este trabajo se origina en un proyecto de investigación sobre la Medición y el Análisis de la Distribución del Ingreso en los países de América Latina, que es realizado conjuntamente por la Comisión Económica para América Latina y el Development Research Center del Banco Mundial.



ANALISIS DE DESCOMPOSICION: UNA GENERALIZACION DEL METODO DE THEIL

I. Introducción

En el contexto de la teoría de la información H. Theil (1967, 1972) define el concepto de entropía o información esperada de una distribución de probabilidades como el valor esperado del logaritmo de las probabilidades con signo negativo.

Este concepto, que corresponde básicamente a una medida de incertidumbre o desorden, tuvo su origen en la ciencia física, pero ha tenido aplicaciones en el campo de la Economía y Política. Dos ejemplos de estas aplicaciones son el Índice de Concentración Industrial de Hirschman (1945) y de Herfindahl (1950) y el Índice de Cohesión de Rice (1928).

A partir de este concepto de entropía, Theil (1967, 1972) deriva una medida de la desigualdad de los ingresos de una determinada población. Esta medida se define como la información esperada del mensaje que transforma porcentajes poblacionales en participaciones de ingresos. Esta medida, conocida como el Índice de Theil, tiene ciertas propiedades de descomposición que la hacen particularmente atractiva para el análisis multivariado de la desigualdad de los ingresos mediante una evaluación de la influencia o efecto "atribuible" a diferentes factores en la generación de esa desigualdad.^{1/}

Si bien este análisis de descomposición de datos no llega a aventurarse en el campo de los mecanismos de causalidad entre el ingreso y las distintas características o variables "explicativas", él facilita e ilumina la formulación y verificación de hipótesis con respecto a estas relaciones causales.

^{1/} Estas propiedades de descomposición se deben a la aditividad de la función de probabilidad de la información propuesta por Shannon (1948) como $h(p) = -\log p$

La aplicación del Índice de Theil, como método de análisis de descomposición, ha sido ya utilizada en algunos análisis empíricos de distribución de ingresos. (Fishlow, 1972; Van Ginneken, 1975; Ullman Chiswick, 1976). Sin embargo, estas aplicaciones han sido condicionadas por las limitaciones y características de las bases de datos utilizados.

En el presente trabajo se propone una generalización formal del método de descomposición de Theil y se avanza en la determinación de sus propiedades y en la interpretación de las interacciones, con el objeto de habilitar su aplicación a diversas bases de datos sobre distribución de ingresos.

II. Formulación de Theil

El coeficiente de entropía de la distribución de ingresos de una determinada población está dado por

$$(1) \quad H(y) = - \sum_{u=1}^N y_u \text{Log } y_u$$

en que N es el número de individuos en la población e y_u la participación del individuo u-ésimo en el ingreso total.

Este coeficiente es una medida de desigualdad que fluctúa entre 0 y Log N, extremos que corresponden a los casos de perfecta desigualdad y perfecta igualdad, respectivamente.^{2/}

Theil (1967) transforma este coeficiente de entropía en una medida de desigualdad restando su valor de su propio valor máximo.

$$(2) \quad T = \text{Log } N - H(y) = \sum_{u=1}^N y_u \text{Log } \frac{y_u}{n_u} = \sum_{u=1}^N y_u \text{Log } \frac{y_u}{1/N}$$

expresión conocida como índice de desigualdad de los ingresos de Theil, en que n_u representa la participación de la u-ésima unidad en la población total y corresponde a $1/N$ para cada individuo. Este índice

^{2/} Véase el Apéndice Matemático 2, para una demostración de estos límites.

de desigualdad está acotado en ambos sentidos: fluctúa entre cero para el caso de perfecta igualdad ($y_u = 1/N$ para todo u) y $\text{Log } N$ para el caso de perfecta desigualdad ($y_u = 1$ para $u=k$ e $y_u = 0$ para todo $u \neq k$).^{3/} El hecho de que el campo de variación de este índice no sea invariante con respecto al tamaño de la población representa un inconveniente para la comparación del grado de desigualdad en la distribución del ingreso entre poblaciones de distinto tamaño. Una forma de obviar este inconveniente es estandarizar el índice dividiendo el valor que éste tome por el logaritmo natural del tamaño de la población respectiva

$$(3) \quad T^* = \frac{T}{\text{Log } N}$$

de manera que el valor efectivo de cada índice de Theil, quede expresado como proporción de su propio valor máximo.

^{3/} Si la población N se particiona en G grupos excluyentes entre sí S_g dentro de los cuales no existe o no se considera la variabilidad en el ingreso, entonces el campo de variación del índice de Theil disminuye. Si todos los grupos son de igual tamaño ($n_g = \frac{N}{G}$) el campo de variación será $0 \leq T \leq \text{Log } G \leq \text{Log } N$. Si los grupos son de distinto tamaño, entonces el campo de variación será $0 \leq T \leq \text{Log } M_g \leq \text{Log } N$ en que $M_g = \text{Max} \left\{ \frac{N}{N_g} \right\}$

Esta disminución del campo de variación impuesta por la existencia de grupos dentro de los cuales no se considera la variabilidad de ingreso es particularmente importante debido a que en la mayoría de los estudios empíricos, la población se obtiene a partir de la expansión de una muestra aplicando a cada observación muestral su respectivo coeficiente de expansión y ello genera, por lo tanto, grupos en la población dentro de los cuales no existe variación en el ingreso.

Una de las ventajas del índice de Theil reside en sus propiedades de descomposición. Theil (1967, 1972) demuestra que si la población se particiona en G grupos ($S_1, S_2 \dots S_g$) de manera que cada uno de los individuos que la constituyen pertenezca a uno y sólo uno de estos grupos, el índice de Theil puede desagregarse en dos componentes:

i) un componente (B) que representa el aporte a la desigualdad total de la desigualdad entre los promedios de ingresos de los distintos grupos,

ii) Otro componente (W) que representa el aporte a la desigualdad total de las desigualdades en el interior de cada uno de los grupos. Puede descomponerse en

$$(4) \quad T = \sum_{g=1}^G y_g \text{Log} \frac{y_g}{n_g} + \sum_{g=1}^G y_g \sum_{u \in S_g} \frac{y_{gu}}{y_g} \text{Log} \frac{\frac{y_{gu}}{y_g}}{1/N_g}$$

$$(5) \quad T = B_g + W_g \quad \text{en que}$$

$$(6) \quad B_g = \sum_{g=1}^G y_g \text{Log} \frac{y_g}{n_g} \quad y$$

$$(7) \quad W_g = \sum_{g=1}^G y_g \sum_{u \in S_g} \frac{y_{gu}}{y_g} \text{Log} \frac{\frac{y_{gu}}{y_g}}{1/N_g}$$

En que y_{gu} es la participación de la u -ésima observación del grupo g en el ingreso total e y_g y n_g las participaciones del grupo g en el ingreso total y la población total respectivamente

$$\left(y_g = \sum_{u \in S_g} y_{gu}, \quad n_g = \sum_{u \in S_g} n_u \right)$$

La expresión (4) puede reescribirse como

$$(8) \quad T = B_g + \sum_{g=1}^G Y_g T_g \quad \text{en que}$$

$$(9) \quad T_g = \sum_{u \in S_g} \frac{y_{gu}}{y_g} \text{Log} \frac{\frac{y_{gu}}{y_g}}{1/N_g}$$

El primer término de las expresiones (4), (5) y (8) corresponde a la componente "entre grupos" (Theil, 1967, 1972) o "parte explicada" de la desigualdad total T (van Ginneken, 1975). El segundo término de estas expresiones corresponde a la componente "dentro de grupos" (Theil, 1967, 1972) o "parte no explicada" de la desigualdad total por la partición G (van Ginneken, 1975). Este último término corresponde a un promedio ponderado de medidas de desigualdad de Theil computadas dentro de cada uno de los grupos, en que las ponderaciones corresponden a las participaciones de los grupos en el ingreso total.

Es importante recordar que el campo de variación de T_g es función del tamaño de cada grupo, ya que T_g varía entre cero y $\text{Log } N_g$ (siendo N_g el tamaño del grupo "g"). Por lo tanto, en la medida en que los grupos difieran en tamaño, diferirá también el rango de variación del índice respectivo. Esto dificulta la utilización directa de los índices de Theil calculados dentro de cada uno de los grupos para comparar el grado de desigualdad de sus distribuciones de ingreso y hace necesario el proceso de estandarización mencionado anteriormente.

Esta propiedad de descomposición del índice de Theil hace extremadamente atractiva su aplicación al análisis de la desigualdad de los ingresos. Se puede particionar el universo muestral de acuerdo con una o más variables clasificatorias - graduadas o no - y la descomposición de la desigualdad total permite derivar medidas de la contribución o efecto sobre la desigualdad de ingresos atribuible a cada una de las variables utilizadas para particionar la población.

III. El caso particular de dos características clasificatorias

Si la partición de la población se hace de acuerdo con una característica estratificadora i que puede tomar \bar{i} valores distintos generando por tanto \bar{i} grupos distintos, S_i , con N_i individuos cada uno, entonces, siguiendo la expresión (4), el índice de Theil se puede escribir como:

$$(10) \quad T = \sum_{i=1}^{\bar{i}} y_i \text{ Log } \frac{y_i}{n_i} + \sum_{i=1}^{\bar{i}} y_i \sum_{u \in S_i} \frac{y_{iu}}{y_i} \text{ Log } \frac{y_{iu} y_i}{1/N_i}$$

$$T = B_i + W_i$$

$$(11) \quad W_i = \sum_{i=1}^{\bar{i}} y_i T_i$$

En que y_i y n_i corresponden a las participaciones del grupo i en el ingreso total y en la población total, respectivamente, y T_i al índice de Theil calculado dentro del grupo i .

B_i representa aquella parte de la desigualdad total "explicada" por la variable i (componente entre grupos S_i). W_i representa aquella parte de la desigualdad total "no explicada" por la variable i (componente dentro de los grupos S_i) y corresponde a un promedio ponderado de las medidas de desigualdad computadas dentro de cada grupo i .

$\frac{B_i}{T}$ y $\frac{W_i}{T}$ indican las proporciones de la desigualdad total "explicada" y "no explicada", respectivamente, por la variable i .^{4/} Si se introduce una segunda variable

^{4/} El uso de la terminología "explicada" y "no explicada" tomada de van Ginneken (1975) no implica necesariamente la existencia de una relación de causalidad entre el ingreso y la variable i .

estratificadora j que puede tomar \bar{j} valores distintos, entonces - por analogía con el caso anterior - el índice de Theil se puede escribir como

$$(12) \quad T = B_j + W_j \quad \text{en que los términos se definen en forma análoga al caso anterior.}$$

a) Contribuciones a la desigualdad total

Se define B_i y B_j como las contribuciones individuales brutas de las variables i y j a la desigualdad total cuando cada una de ellas es independientemente considerada como la variable clasificatoria.

Si se clasifica la población de acuerdo con ambas características simultáneamente se puede obtener la contribución conjunta de ambas variables. Esta clasificación generará $\bar{i} \cdot \bar{j}$ grupos distintos S_{ij} con N_{ij} individuos cada uno. En este caso el índice de Theil se puede escribir como

$$(13) \quad T = \sum_{i=1}^{\bar{i}} \sum_{j=1}^{\bar{j}} y_{ij} \text{Log} \frac{y_{ij}}{n_{ij}} + \sum_{i=1}^{\bar{i}} \sum_{j=1}^{\bar{j}} y_{ij} \sum_{u \in S_{ij}} \frac{y_{iju}}{y_{ij}} \text{Log} \frac{y_{iju} / y_{ij}}{1/N_{ij}}$$

$$T = B_{ij} + W_{ij}$$

$$(14) \quad W_{ij} = \sum_i \sum_j y_{ij} T_{ij}$$

En que y_{ij} y n_{ij} son las participaciones del grupo S_{ij} en el ingreso y la población respectivamente, y T_{ij} el valor del índice de Theil para cada grupo S_{ij}

Se define B_{ij} como la contribución conjunta de las variables i y j y corresponde a aquella parte de la desigualdad total "explicada" en forma conjunta por las variables i y j (componente entre grupos S_{ij}). W_{ij} representa aquella parte de la desigualdad total "no explicada" por las variables i y j (componente dentro de los grupos S_{ij}) y corresponde a un promedio ponderado de los niveles de desigualdad dentro de cada uno de los grupos S_{ij} , en que las ponderaciones son sus respectivas participaciones en el ingreso total.

Sin embargo, el índice de Theil también puede descomponerse en este caso de las siguientes dos maneras

$$(15) \quad T = B_i + B_j^i + W_{ij} \quad o$$

$$(16) \quad T = B_j + B_i^j + W_{ij}$$

En que:

$$(17) \quad B_i^j = \sum_j y_j \sum_i \frac{y_{ij}}{y_j} \text{Log} \frac{y_{ij} y_j}{n_{ij} n_j}$$

$$(18) \quad B_j^i = \sum_i y_i \sum_j \frac{y_{ij}}{y_i} \text{Log} \frac{y_{ij} y_i}{n_{ij} n_i}$$

B_i^j y B_j^i se definen como las Contribuciones Marginales de la variable i dada la variable j , y de la variable j dada la variable i respectivamente. En las expresiones (17) y (18) se puede observar que B_i^j (B_j^i) corresponde a un promedio ponderado de la desigualdad entre los distintos subgrupos definidos por la variable i (j) para cada uno de los grupos definidos por la variable j (i), en que las ponderaciones corresponden a las participaciones en el ingreso total de cada uno de los grupos S_j (S_i).

A partir de las expresiones (13), (15) y (16) vemos que

$$(19) \quad B_{ij} = B_i + B_j^i = B_j + B_i^j$$

es decir, la contribución conjunta (B_{ij}) de las variables i y j puede expresarse siempre como la contribución individual de una de ellas (B_i o B_j) más la contribución marginal de la otra dada la primera (B_j^i o B_i^j)

A partir de las siete últimas expresiones se puede observar que la suma de las contribuciones individuales brutas de cada una de las variables tomadas en forma independiente ($B_i + B_j$) no es necesariamente igual a la contribución conjunta de ambas variables (B_{ij}), lo que equivale a decir que la contribución individual bruta de una variable (B_i o B_j) no es necesariamente igual a su contribución marginal dada la otra (B_j^i o B_i^j). En otras palabras, si las variables interactúan, la contribución conjunta diferirá de la suma de las contribuciones individuales brutas y cada contribución individual bruta diferirá de la contribución marginal respectiva.

b) Interacciones entre variables

Se define el coeficiente de interacción entre las dos variables (I_{ij}) como la diferencia entre la parte de la desigualdad total explicada en forma conjunta por ambas variables y la suma de las partes explicadas por cada una de ellas consideradas en forma independiente

$$(20) \quad I_{ij} = B_{ij} - B_i - B_j \quad \text{a partir de (19) vemos que}$$

$$(21) \quad I_{ij} = B_i^j - B_i \equiv B_j^i - B_j$$

Este coeficiente constituye la contribución de la interacción a la desigualdad total.

Existen dos condiciones suficientes pero no necesarias bajo las cuales estas dos variables tienen una interacción nula y por lo tanto sus contribuciones son independientes.

i) Que la distribución de la población de acuerdo con la característica i sea independiente de la distribución de la población de acuerdo con la característica j . Es decir, que las probabilidades condicionadas de i dado j , y de j dado i , sean iguales a las probabilidades marginales de i y de j , respectivamente, ya que esto implica que $n_{ij} = n_i \cdot n_j$

ii) Que la participación en el ingreso total del grupo S_{ij} sea igual a la participación del grupo S_i multiplicado por la participación del grupo S_j : $y_{ij} = y_i \cdot y_j$

Si se dan ambas condiciones simultáneamente entonces

$$(22) \quad B_{ij} = \sum_i \sum_j y_{ij} \text{ Log } \frac{y_{ij}}{n_{ij}} = \sum_i y_i \text{ Log } \frac{y_i}{n_i} + \sum_j y_j \text{ Log } \frac{y_j}{n_j}$$

Es decir, la contribución conjunta de ambas variables sería igual a la suma de sus contribuciones individuales brutas y las contribuciones marginales de ambas variables iguales a sus contribuciones individuales brutas. Dado (20) y (21), esto implica que la interacción entre ellos (I_{ij}) es nula. Por lo tanto, las dos condiciones mencionadas anteriormente implican que las contribuciones de ambas variables son mutuamente independientes. Sin embargo, una interacción nula entre ellas no implica necesariamente que se den las dos condiciones antes señaladas. 5/

5/ Esto debido a que $n_{ij} = n_i \cdot n_j \wedge y_{ij} = y_i \cdot y_j \Rightarrow$

$$\frac{y_{ij}}{n_{ij}} = \frac{y_i}{n_i} \cdot \frac{y_j}{n_j} \quad \text{pero}$$

$$\frac{y_{ij}}{n_{ij}} = \frac{y_i}{n_i} \cdot \frac{y_i}{n_j} \not\Rightarrow n_{ij} = n_i \cdot n_j \wedge y_{ij} = y_i \cdot y_j$$

La expresión (20) implica que la contribución conjunta de ambas variables es igual a la suma de sus contribuciones individuales brutas más la interacción entre ambas. Por lo tanto, será mayor, igual o menor que esta suma según que la interacción sea mayor, igual o menor que cero, respectivamente. Similarmente, a partir de la expresión (21) vemos que la contribución marginal de cada una de las variables es igual a su contribución individual bruta más la interacción. La contribución marginal será, por lo tanto, mayor, igual o menor que la contribución individual, dependiendo del signo de la interacción.

A partir de (20) el término de interacción se puede descomponer en dos componentes. Uno de ellos está relacionado con las participaciones de los diferentes grupos en el ingreso total y el otro con las participaciones de los diferentes grupos en la población total

$$\begin{aligned}
 (23) \quad I_{ij} &= B_{ij} - B_i - B_j = \sum_i \sum_j y_{ij} \operatorname{Log} \frac{y_{ij}}{n_{ij}} - \\
 &\sum_i y_i \operatorname{Log} \frac{y_i}{n_i} - \sum_j y_j \operatorname{Log} \frac{y_j}{n_j} \\
 I_{ij} &= \sum_i \sum_j y_{ij} \operatorname{Log} \left(\frac{y_i}{n_i} \cdot \frac{y_j}{n_j} \cdot \frac{y_{ij}}{y_i \cdot y_j} \cdot \frac{n_i n_j}{n_{ij}} \right) - \\
 &\sum_i y_i \operatorname{Log} \frac{y_i}{n_i} - \sum_j y_j \operatorname{Log} \frac{y_j}{n_j} \\
 (24) \quad I_{ij} &= \sum_i \sum_j y_{ij} \left[\operatorname{Log} \frac{y_{ij}}{y_i \cdot y_j} - \operatorname{Log} \frac{n_{ij}}{n_i \cdot n_j} \right]
 \end{aligned}$$

El primer término de la expresión (21) no puede tomar valores negativos.^{6/} El segundo término de esta expresión puede tomar, en cambio, valores positivos, nulos o negativos. Si la distribución de la población de acuerdo con la variable i es independiente de la distribución de la población de acuerdo con la variable j , es decir, si no hay asociación estadística entre ambas variables, entonces $n_{ij} = n_i \cdot n_j$ y por lo tanto el segundo término de la expresión (24) se hace nulo quedando la interacción reducida a

$$(25) \quad I_{ij} = \sum_i \sum_j y_{ij} \text{ Log } \frac{y_{ij}}{y_i \cdot y_j} \geq 0$$

que es necesariamente una expresión no negativa. Por lo tanto, si las dos variables son estadísticamente independientes la contribución conjunta de ambas será siempre mayor o igual que la suma de sus contribuciones individuales brutas y la contribución marginal de cada una de ellas será siempre mayor o igual que la respectiva contribución individual bruta. Sin embargo, si las dos variables están estadísticamente correlacionadas, entonces, la interacción puede tomar valores positivos, negativos o nulos y entonces la contribución conjunta puede exceder, o ser excedida, o ser igual a la suma de las contribuciones individuales brutas.

Con respecto al campo de variación del coeficiente de interacción, un caso extremo se presenta cuando existe una perfecta asociación o correlación entre las dos variables y por tanto, la contribución conjunta a la desigualdad es igual a cada una de las contribuciones individuales brutas. A partir de (20) vemos que en este caso extremo $I_{ij} = -B_{ij}$, éste constituye el límite inferior para el coeficiente de interacción. A partir de la misma expresión se comprueba que el máximo valor que puede tomar el coeficiente de interacción es B_{ij} y ello ocurre cuando las contribuciones individuales de ambas variables son nulas. Por lo tanto se puede concluir que la interacción entre ambas variables está acotada por

$$(26) \quad - \text{Log } N \leq -T \leq -B_{ij} \leq I_{ij} \leq B_{ij} \leq T \leq \text{Log } N$$

^{6/} Véase el Apéndice Matemático N^o 3, para una demostración de este punto.

Los cuadros 1 y 2 proveen ejemplos numéricos de las dos situaciones extremas descritas.

La interacción también se puede presentar como proporción de la contribución conjunta; como ésta es el valor máximo de la interacción, resulta un coeficiente estandarizado, útil para efectuar comparaciones entre diferentes bases de datos:

$$(27) \quad I_{ij}^* = \frac{I_{ij}}{B_{ij}}$$

c) Interacciones negativas y positivas

El significado de las interacciones negativas y positivas puede ser aprehendido mediante la explicación de algunos casos extremos.

Considérese el caso en que sólo la variable i tenga efectivamente un impacto causal en la determinación del nivel de ingresos, en tanto que la variable j no tiene ningún efecto autónomo sobre él. Si ambas variables se hallan estadísticamente asociadas, las distribuciones de la población y de los ingresos de acuerdo con la variable i no serán independientes de las correspondientes distribuciones de acuerdo con la variable j . En este caso, al clasificar la población de acuerdo con la variable j , B_j será positivo, aun cuando esta variable no ejerza un impacto real sobre los ingresos, debido a que ella se "apropia" o "captura" parte de la influencia de i sobre el ingreso, a través de su asociación con esta variable. En este caso, la suma de las contribuciones individuales ($B_i + B_j$) incluye dos veces aquella parte de la influencia de i capturada por j a través de su asociación con i ; esa suma excede por lo tanto, la contribución conjunta de ambas variables, dando lugar a una interacción negativa. En este caso, que se ilustra con el ejemplo numérico del cuadro 3, la contribución marginal de j controlando por i es nula, como resultado de que j no tiene influencia por sí misma sobre el ingreso y no contribuye adicionalmente a explicar las desigualdades. Esto da origen a una interacción negativa que neutraliza totalmente la contribución individual de j :

$$(28) \quad I_{ij} = B_j^i - B_j = - B_j < 0$$

$$(29) \quad B_{ij} = B_i + B_j + I_{ij} = B_i$$

El ejemplo numérico del cuadro 4 ilustra una situación similar, aunque no extrema, en que ambas variables se encuentran asociadas y la influencia de i sobre el ingreso es considerablemente mayor que la que ejerce j , aunque esta variable tenga alguna influencia por sí misma. También en este caso la variable j captura parte de la influencia de i en la medida en que se halla asociada con ella.^{7/} La suma $(B_i + B_j)$ incluye, por lo tanto, duplicaciones que no están presentes en la contribución conjunta B_{ij} y que están representadas por una interacción negativa. El cuadro 1 ilustra, finalmente, el caso extremo en que las variables i y j están perfectamente asociadas y la desigualdad resulta, por lo tanto, totalmente explicada por cualquiera de ellas; esto da lugar a una interacción negativa que es, en consecuencia, igual a la contribución conjunta. En este caso se puede considerar que ambas variables constituyen, en realidad, una sola variable explicativa.

La presencia de interacciones negativas revela, en todos los casos, una asociación subyacente entre las variables, que posibilita que la desigualdad sea, en alguna medida, explicada por una de las variables en representación de las asociadas.

La posibilidad de interacciones positivas es ilustrada por los ejemplos numéricos de los cuadros 2 y 5.

Considérese el caso extremo en que las distribuciones de acuerdo con cada una de las variables no registren desigualdades; éstas sólo pueden ser explicadas por determinadas combinaciones de las variables, es decir, por la interacción entre ambas (cuadro 2). En este caso, la interacción es positiva e igual a la contribución conjunta de las

^{7/} También la variable i captura, en este caso, parte la escasa influencia de j .

variables, mientras que sus contribuciones individuales brutas son nulas. Este mismo resultado puede darse aun cuando exista asociación entre las variables (en la diagonal del cuadro 2).

Puede concebirse, finalmente, un caso en que cada variable registre desigualdades, cuando se clasifica a la población sólo con respecto a ella, pero esas desigualdades se originen en una combinación particular de ambas variables (cuadro 5). En este caso, la contribución conjunta excede la suma de las contribuciones individuales, en tanto existe una interacción positiva, pues sólo cuando se clasifica la población de acuerdo con ambas variables simultáneamente se logra captar la interacción existente entre ellas.

La presencia de interacciones positivas debe interpretarse como que la influencia de una variable sobre el ingreso no es independiente del valor que tome la otra variable, ya sea que esa influencia se ejerza sólo para determinados valores de la otra variable o que sea mayor cuando ésta toma esos valores.^{8/} Esta situación implica que la combinación de las variables agrega poder explicativo que no es captado por ninguna de ellas en forma individual.

En síntesis, el coeficiente de interacción (I_{ij}) obtenido a partir del análisis de descomposición de Theil, es el resultado neto de dos tipos de factores de distinta naturaleza. El primero, la asociación estadística entre las variables que puede generar contribuciones individuales brutas que exceden a las correspondientes contribuciones marginales, dando lugar a interacciones negativas entre ellas. El segundo, presenta cuando parte de la determinación del nivel de ingreso no puede explicarse en base a las distintas variables tomadas en forma independiente sino que es el resultado de combinaciones de ellas, puede dar origen a interacciones positivas.

Dependiendo de cuál de estas dos fuentes de interacciones predomine será el signo del coeficiente de interacción obtenido del análisis de descomposición. La presencia de interacciones positivas indica inequívocamente la existencia de la segunda fuente de interacciones. La posibilidad de detectar este tipo de causalidad constituye una innegable ventaja del método de descomposición de Theil.

^{8/} Un ejemplo de esto puede encontrarse entre la educación y el sector (formal-informal) o el área, cuando la influencia de la educación sobre el ingreso es mayor en el sector formal que en el informal, o es mayor en las áreas urbanas que en las rurales.

Cuadro 1

EJEMPLO NUMERICO DEL CASO EXTREMO DE PERFECTA ASOCIACION ENTRE LAS VARIABLES EXPLICATIVAS DEL INGRESO

$j \backslash i$	$i = 1$	$i = 2$	$i = 3$	N_j	Y_j
$j = 1$	1,1	0,0	0,0	1	1
$j = 2$	0,0	1,2	0,0	1	2
$j = 3$	0,0	0,0	1,3	1	3
N_i	1	1	1	$N = 3$	
Y_i	1	2	3	$Y = 6$	

$$B_i = B_j = 0,200804$$

$$B_i + B_j = 0,401608$$

$$B_{ij} = 0,200804$$

$$I_{ij} = -0,200804 = -B_{ij} = -B_i = -B_j$$

$$B_i^j = B_j^i = 0$$

Cuadro 2

EJEMPLO NUMERICO DEL CASO EXTREMO EN QUE SOLO LAS VARIABLES
COMBINADAS EXPLICAN LA DESIGUALDAD.

$j \backslash i$	$i = 1$	$i = 2$	$i = 3$	N_j	Y_j
$j = 1$	1,3	1,1	1,1	3	5
$j = 2$	1,1	1,3	1,1	3	5
$j = 3$	1,1	1,1	1,3	3	5
N_i	3	3	3	$N = 9$	
Y_i	5	5	5	$Y = 15$	

$$B_i = B_j = 0$$

$$B_i = B_j = 0$$

$$B_{ij} = 0,1483$$

$$I_{ij} = 0,1483 = B_{ij}$$

Cuadro 3

EJEMPLO NUMERICO DEL CASO EXTREMO EN QUE LAS VARIABLES SE HALLAN ASOCIADAS PERO SOLO UNA DE ELLAS TIENE INFLUENCIA SOBRE EL INGRESO

$j \backslash i$	$i = 1$	$i = 2$	$i = 3$	N_j	Y_j
$j = 1$	3,3	1,2	1,2	5	8
$j = 2$	1,1	3,6	1,3	5	10
$j = 3$	1,1	1,2	3,9	5	12
N_i	5	5	5	$N = 15$	
Y_i	5	10	15	$Y = 30$	

$$B_i = 0,087208 \quad B_i^j = 0,073784$$

$$B_j = 0,013424 \quad B_j^i = 0$$

$$B_{ij} = 0,087208$$

$$I_{ij} = 0,013424$$

Cuadro 4

EJEMPLO NUMERICO DE UN CASO EN QUE LAS VARIABLES EXPLICATIVAS SE ENCUENTRAN ASOCIADAS Y UNA DE ELLAS TIENE MAYOR INFLUENCIA SOBRE EL INGRESO a/

$j \backslash i$	$i = 1$	$i = 2$	$i = 3$	N_j	Y_j
$j = 1$	3,12	1,7	1,10	5	29
$j = 2$	1,5	3,24	1,11	5	40
$j = 3$	1,6	1,9	3,36	5	51
N_i	5	5	5	$N = 15$	
Y_i	23	40	57	$Y = 120$	

$$B_i = 0,062166$$

$$B_j = 0,025536$$

$$B_i + B_j = 0,087702$$

$$B_{ij} = 0,067189$$

$$I_{ij} = 0,020514$$

$$B_j^i = 0,05022$$

$$B_i^j = 0,041652$$

a/ La relación funcional supuesta para construir el ejemplo es

$$Y = 3i + j$$

Cuadro 5

EJEMPLO NUMERICO DEL CASO EXTREMO EN QUE SOLO UNA DETERMINADA
COMBINACION DE LAS VARIABLES EXPLICA
TODA LA DESIGUALDAD

$j \backslash i$	$i = 1$	$i = 2$	$i = 3$	N_j	Y_j
$j = 1$	1,10	1,1	1,1	3	12
$j = 2$	1,1	1,1	1,1	3	3
$j = 3$	1,1	1,1	1,1	3	3
N_i	3	3	3	3	$N = 9$
Y_i	12	3	3	3	$Y = 18$

$$B_i = B_j = 0,231048$$

$$B_i + B_j = 0,462097$$

$$B_{ij} = 0,586068$$

$$I_{ij} = 0,123971$$

IV. El caso general de R características estratificadoras

Todos los conceptos enunciados anteriormente pueden generalizarse para el caso de R variables C_r ($r = 1, 2, \dots, R$) en que cada una de ellas puede tomar \bar{C}_i valores distintos. Por lo tanto, al clasificar la población de acuerdo con las R variables simultáneamente se generan

$$\prod_{r=1}^R \bar{C}_r = \bar{C}_1 \cdot \bar{C}_2 \cdot \dots \cdot \bar{C}_R \quad \text{celdas o grupos diferentes}$$

$$S_{C_1, C_2, \dots, C_R} \quad \text{con} \quad N_{C_1, C_2, \dots, C_R} \quad \text{individuos cada uno}$$

En estas circunstancias, el índice de Theil se puede escribir como

$$(30) \quad T = \sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2, \dots, C_R} \text{Log} \frac{y_{C_1, C_2, \dots, C_R}}{n_{C_1, C_2, \dots, C_R}} +$$

$$\sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2, \dots, C_R} \sum_{u \in S_{C_1, C_2, \dots, C_R}} \frac{y_{C_1, C_2, \dots, C_R, u}}{y_{C_1, C_2, \dots, C_R}}$$

$$\text{Log} \frac{y_{C_1, C_2, \dots, C_R, u}}{y_{C_1, C_2, \dots, C_R}} \cdot \frac{1}{N_{C_1, C_2, \dots, C_R}}$$

$$T = B_{1, 2, \dots, R} + W_{1, 2, \dots, R}$$

$$(30) \quad W_{1, 2, \dots, R} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2, \dots, C_R} T_{C_1, C_2, \dots, C_R}$$

En que todos los términos quedan definidos por analogía al caso particular de dos variables analizado en la sección anterior.

a) Contribuciones individuales brutas conjuntas y marginales de las variables

En este caso general, el índice de desigualdad de Theil, también puede descomponerse en la suma de (R+1) términos de R! maneras distintas.^{9/} Una de estas R! maneras es la siguiente

$$(32) \quad T = B_1 + B_2^1 + B_3^{1,2} \dots + B_q^{1,2\dots(q-1)} + \dots + B_R^{1,2\dots(R-1)} + W_{1,2\dots R}$$

Las restantes (R! - 1) maneras de descomponer este índice se obtienen variando el orden de las variables.^{10/} Dado que el análisis que sigue es exactamente análogo para cualquiera de las R! posibles descomposiciones del índice T, escogeremos por conveniencia la ordenación presentada en (32). Nuevamente las definiciones de los distintos términos de (32) son análogas a las del caso particular de dos variables ya analizado. Por lo tanto, nos limitaremos a definir el primer término, el término genérico de orden q y el último término o residuo de la expresión (27).

$$(33) \quad B_1 = \sum_{C_1=1}^{\bar{C}_1} y_{C_1} \text{Log} \frac{y_{C_1}}{n_{C_1}}$$

B_1 representa la componente de desigualdad entre los promedios de ingresos de los \bar{C}_1 grupos definidos por la variable C_1 y corresponde a la contribución individual bruta de esta variable a la desigualdad total. B_2^1 representa la contribución marginal de la segunda variable dado que el efecto de la primera ya ha sido considerado. Por lo tanto,

^{9/} La demostración matemática de esta descomposición aparece en el apéndice matemático 3.

^{10/} Otras de las R! alternativas de descomposición son las siguientes

i) $T = B_2 + B_1^2 + B_3^{1,2} + \dots + B_R^{1,2\dots(R-1)}$

ii) $T = B_1 + B_3^1 + B_2^{1,3} + \dots + B_R^{1,2\dots(R-1)}$

iii) $T = B_2 + B_3^2 + B_1^{2,3} + \dots + B_R^{1,2\dots(R-1)}$

la suma de B_1 y B_2^1 corresponde a la contribución conjunta de ambas variables. Análogamente, $B_3^{1,2}$ representa la contribución marginal de la tercera variable dadas las dos primeras. Por lo tanto, la suma de B_1 , B_2^1 y $B_3^{1,2}$ corresponde a la contribución conjunta de las tres variables. En términos generales, si cada logaritmo involucrado está definido, el término genérico de orden q ($1 \leq q \leq R$) se define como

$$(34) \quad B_q^{1,2,\dots,(q-1)} = \sum_{C_1} \sum_{C_2 \dots C_{(q-1)}} y_{C_1 C_2 \dots C_{(q-1)}} \sum_{C_q} \frac{y_{C_1, C_2 \dots C_q}}{y_{C_1, C_2 \dots C_{(q-1)}}} \text{Log} \frac{\frac{y_{C_1, C_2 \dots C_q}}{n_{C_1, C_2 \dots C_q}}}{\frac{y_{C_1, C_2 \dots C_{(q-1)}}}{n_{C_1, C_2 \dots C_{(q-1)}}}}$$

Este término de orden q representa la contribución marginal de la q -ésima variable controlando por las $(q-1)$ variables anteriores. La suma de la contribución individual bruta de la primera variable, más la contribución marginal de la segunda dada la primera, más la contribución marginal de la tercera dada las dos primeras, y así sucesivamente hasta agregar la contribución marginal de la q -ésima variable dadas las $(q-1)$ anteriores corresponde a la contribución conjunta de las q variables.

$$(35) \quad B_{1,2,\dots,q} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_q} y_{C_1, C_2 \dots C_q} \text{Log} \frac{y_{C_1, C_2 \dots C_q}}{n_{C_1, C_2 \dots C_q}} =$$

$$B_1 + B_2^1 + B_3^{1,2} + \dots + B_q^{1,2,\dots,(q-1)} \quad \underline{11/}$$

11/ El término $B_{1,2,\dots,q}$ también puede escribirse de $q!$ maneras distintas. Una de ellas es la presentada en la expresión (30), las restantes se obtienen cambiando el orden de las variables.

Finalmente, el residuo de la expresión (32), es decir $W_{1,2..R}$, representa aquella parte de la desigualdad total no explicada por las R variables utilizadas y se define como

$$(36) \quad W_{1,2..R} = \sum_{C_1} \sum_{C_2..C_R} y_{C_1,C_2..C_R} \sum_{u \in S_{C_1,C_2..C_R}} \frac{y_{C_1,C_2..C_R,u}}{y_{C_1,C_2..C_R}}$$

$$\text{Log} \frac{\frac{y_{C_1,C_2..C_R,u}}{y_{C_1,C_2..C_R}}}{1/N_{C_1,C_2..C_R}}$$

$$(37) \quad W_{1,2..R} = \sum_{C_1} \sum_{C_2..C_R} y_{C_1,C_2..C_R} T_{1,2..R}$$

Es decir, la parte de la desigualdad total "no explicada" por las R variables utilizadas en forma conjunta corresponde a un promedio ponderado de los niveles de desigualdad al interior de cada uno de los grupos $S_{C_1,C_2..C_R}$ en que las ponderaciones son las participaciones de cada uno de estos grupos en el ingreso total.

A partir de la expresión (34) se observa que la contribución marginal de la q-ésima variable dada las (q-1) anteriores corresponde a un promedio ponderado de los índices de Theil correspondientes al grado de desigualdad dentro de cada uno de los grupos definidos por las (q-1) variables iniciales. Por lo tanto, la contribución marginal de una variable no puede ser negativa, lo que implica que la parte explicada de la desigualdad total no puede disminuir con la inclusión de una nueva variable. A partir de esta misma expresión, también se observa que la contribución marginal de una variable depende de cuántas y cuáles de las otras variables están siendo controladas y por lo tanto, una variable tendrá tantas contribuciones marginales como conjuntos de las otras variables se controlen para computarla. Los dos casos extremos están representados por la contribución individual bruta en que no se controla por ninguna variable y la contribución marginal de orden R en que se controla por todas las (R-1) variables restantes.

b) Interacciones entre variables

Nuevamente vemos que la suma de las contribuciones individuales brutas de las R variables no corresponde necesariamente a la contribución conjunta de ellas. Se define el Coefficiente de Interacción Total entre q variables ($\hat{I}_{1,2..q}$) como la diferencia entre la contribución conjunta de las q variables y la suma de sus contribuciones individuales brutas

$$(38) \quad \hat{I}_{1,2..q} = B_{1,2..q} - (B_1 + B_2 + \dots + B_q)$$

Dos condiciones son, en forma simultánea, suficientes pero no necesarias para que éstas tengan una interacción total nula y por lo tanto sus contribuciones sean independientes.

i) Que la distribución de la población de acuerdo con cada una de las variables sea independiente de la distribución de la población de acuerdo con las demás. Es decir, que la probabilidad de que se dé un determinado valor de la variable i sea independiente de los valores tomados por los restantes (q-1) variables. Esto implica que

$$n_{C_1, C_2 \dots C_q} = n_{C_1} \cdot n_{C_2} \cdot \dots \cdot n_{C_q}$$

ii) Que la participación en el ingreso total del grupo

$S_{C_1, C_2 \dots C_q}$ sea igual al producto de las participaciones de los grupos

$$S_{C_1, C_2 \dots C_q} \quad y_{C_1, C_2 \dots C_q} = y_{C_1} \cdot y_{C_2} \cdot \dots \cdot y_{C_q}$$

Si estas dos condiciones se dan simultáneamente entonces

$$(39) \quad B_{1,2..q} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_q} y_{C_1, C_2 \dots C_q} \text{Log} \frac{y_{C_1, C_2 \dots C_q}}{n_{C_1, C_2 \dots C_q}} =$$

$$\sum_{C_1} y_{C_1} \text{Log} \frac{y_{C_1}}{n_{C_1}} + \sum_{C_2} y_{C_2} \text{Log} \frac{y_{C_2}}{n_{C_2}} + \dots + \sum_{C_q} y_{C_q} \text{Log} \frac{y_{C_q}}{n_{C_q}}$$

$$B_{1,2..q} = B_1 + B_2 + \dots + B_q$$

Es decir, la contribución conjunta de las q variables es igual a la suma de sus contribuciones individuales brutas y por lo tanto, la interacción total entre ellos definida en (38) es nula.

A partir de la expresión (38) se puede desagregar la interacción total entre q variables en dos componentes. Uno relacionado con las participaciones de los diferentes grupos en el ingreso total y el otro con las participaciones en la población

$$(40) \quad \hat{I}_{1,2..q} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_q} y_{C_1, C_2 \dots C_q} \text{Log} \left(\frac{y_{C_1}}{n_{C_1}} \cdot \frac{y_{C_2}}{n_{C_2}} \dots \frac{y_{C_q}}{n_{C_q}} \cdot \frac{y_{C_1, C_2 \dots C_q}}{y_{C_1} \cdot y_{C_2} \dots y_{C_q}} \cdot \frac{n_{C_1} \cdot n_{C_2} \dots n_{C_q}}{n_{C_1, C_2 \dots C_q}} \right) - \sum_{i=1}^q y_{C_i} \text{Log} \frac{y_{C_i}}{n_{C_i}}$$

$$(41) \quad \hat{I}_{1,2..q} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_q} y_{C_1, C_2 \dots C_q} \left[\text{Log} \frac{y_{C_1, C_2 \dots C_q}}{y_{C_1} \cdot y_{C_2} \dots y_{C_q}} - \text{Log} \frac{n_{C_1, C_2 \dots C_q}}{n_{C_1} \cdot n_{C_2} \dots n_{C_q}} \right]$$

El primer término de la expresión (41) no puede tomar valores negativos. El segundo término de esta expresión puede tomar valores positivos, nulos o negativos. Por lo tanto, si la distribución de la población de acuerdo a cada una de las variables es independiente de la distribución de la población de acuerdo a las demás, el segundo término se hace cero quedando la interacción total reducida al primer término que es número negativo. Esto implica que, si las q variables son estadísticamente independientes, la interacción total entre ellas será siempre mayor o igual que cero. Sin embargo, si las variables no son independientes, esta interacción puede tomar valores positivos, nulos o negativos. Esta interacción total entre q variables puede a su vez descomponerse en la suma de C_2^q interacciones de segundo orden (I_{ij})

correspondientes a todas las combinaciones de dos variables que pueden obtenerse a partir de las q variables, más la suma de C_3^q interacciones de tercer orden (I_{ijk}) correspondientes a todas las combinaciones posibles de tres variables y así sucesivamente hasta incluir la interacción de orden q . Las interacciones de segundo orden entre dos variables se definen como la diferencia entre la contribución conjunta de ellas y la suma de las contribuciones individuales brutas.

$$(42) \quad I_{ij} = B_{ij} - (B_i + B_j)$$

La interacción de tercer orden entre las variables i, j y k se define como la contribución conjunta de ellas menos la suma de sus contribuciones individuales brutas y menos la suma de las interacciones de segundo orden.

$$(43) \quad I_{ijk} = B_{ijk} - (B_i + B_j + B_k) - (I_{ij} + I_{ik} + I_{jk})$$

reemplazando tenemos que

$$(44) \quad I_{ijk} = B_{ijk} - (B_{ij} + B_{ik} + B_{jk}) + (B_i + B_j + B_k)$$

En términos generales, las interacciones de orden P ($2 \leq p \leq q \leq R$) entre p variables (de las cuales habrán C_p^q si se obtienen a partir de q variables) se definen como la contribución conjunta de las p variables menos la suma de sus contribuciones individuales brutas y menos la suma de todas las posibles interacciones de orden menor que p

$$(45) \quad I_{1,2..p} = B_{1,2..p} - (B_1 + B_2 + \dots + B_p) - \sum_{i=1}^{(p-1)} \sum_{j=i+1}^p I_{ij} -$$

$$\sum_{i=1}^{(p-2)} \sum_{j=i+1}^{(p-1)} \sum_{k=j+1}^p I_{ijk} - \dots - \sum_{i=1}^2 \sum_{j=i+1}^3 \dots \sum_{h=k+1}^p \underbrace{I_{ij..h}}_{(p-1) \text{ subíndices}}$$

Reemplazando las interacciones de orden menor que p en términos de contribuciones conjuntas e individuales vemos que la interacción de orden p definida en (45) también puede escribirse como la contribución conjunta de las p variables menos la suma de las C_{p-1}^p contribuciones conjuntas de subconjuntos de $(p-1)$ variables del conjunto inicial de p variables más las C_{p-2}^p contribuciones conjuntas de subconjuntos de $(p-2)$ variables y así sucesivamente hasta el último término que será $(-1)^{p+1}$ por la suma de las contribuciones individuales brutas

$$\begin{aligned}
 (46) \quad I_{1,2..p} &= B_{1,2..p} - \sum_{i=1}^2 \sum_{j=i+1}^3 \dots \sum_{h=k+1}^p \underbrace{B_{ij..h}}_{(p-1)} + \\
 &\sum_{i=1}^3 \sum_{j=i+1}^4 \dots \sum_{h=k+1}^p \underbrace{B_{ij..h}}_{(p-2)} - \dots + (-1)^p \sum_{i=1}^{(p-1)} \sum_{j=i+1}^p B_{ij..} + \\
 &(-1)^{p+1} \sum_{i=1}^p B_i
 \end{aligned}$$

Por lo tanto, la interacción total de orden q definida en (38) como la diferencia entre la contribución conjunta de las q variables y la suma de sus contribuciones individuales brutas es el resultado de una serie de interacciones de segundo orden entre dos variables de tercer orden entre tres variables y así sucesivamente hasta una interacción (no total) de orden q entre las q variables. La interacción total entre las q variables corresponde a la suma algebraica de todas estas interacciones parciales

$$(47) \quad \hat{I}_{1,2..q} = \sum_{i=1}^{(q-1)} \sum_{j=i+1}^q I_{ij} + \sum_{i=1}^{(q-2)} \sum_{j=i+1}^{(q-1)} \sum_{k=j+1}^q I_{ijk} + \dots + I_{1,2..q}$$

Esta interacción total entre las q variables puede ser positiva, nula o negativa dependiendo del resultado neto que originen todas las interacciones parciales que la componen. Es importante destacar que una interacción total nula entre las q variables no implica necesariamente que ellas actúen en forma independiente. Esta interacción total nula entre las q variables es compatible con fuertes interacciones positivas y negativas entre subconjuntos de ellas pero que se cancelan originando una interacción total igual a cero.

Con respecto a las contribuciones marginales, la expresión (44) puede transformarse en

$$(48) \quad B_i^{jk} = B_i + I_{ij} + I_{ik} + I_{ijk}$$

Esto indica que la contribución marginal de la variable i dadas las variables j y k es siempre igual a la contribución individual bruta de i más todas las interacciones posibles entre la variable i y las variables j y k.

Igualmente la expresión (46) puede transformarse en

$$(49) \quad B_i^{j_1 j_2 \dots j_p} = B_i + \sum_{j \neq i} I_{ij} + \sum_{j \neq i} \sum_{k \neq j \neq i} I_{ijk} + \dots + I_{ij_1 j_2 \dots j_p}$$

Lo que significa que la contribución marginal de la variable i dadas las otras (p-1) variables es siempre igual a la contribución individual bruta de i más todas las posibles interacciones de cualquier orden entre la variable i y las demás variables.

Se define la interacción total entre la variable i y las (p-1) variables restantes ($\hat{I}_i^{j_1 j_2 \dots j_p}$) como la suma de todas las posibles interacciones de cualquier orden entre la variable i y las demás (p-1) variables

$$(50) \quad \hat{I}_i^{j_1 j_2 \dots j_p} = \sum_{j \neq i} I_{ij} + \sum_{j \neq i} \sum_{k \neq j \neq i} I_{ijk} + \sum_j \sum_k \sum_l I_{ijkl} + \dots + I_{ij_1 j_2 \dots j_p}$$

Por lo tanto

$$(51) \quad B_i^{jk..p} = B_i + \hat{I}_i^{jk..p}$$

El signo de la interacción total entre la variable i y las demás variables indica si la variable i "en promedio" interactúa positiva o negativamente con las demás variables. La magnitud de esta interacción total refleja la discrepancia entre la contribución marginal y la contribución individual bruta de la variable i .

V. Datos agrupados

a) Introducción

El problema de los datos agrupados se presenta cuando no se dispone de la información respecto de ingresos a un nivel individual sino solamente a un nivel más agregado para grupos de personas.

Sean S_g los grupos de personas para las cuales se dispone de información agregada respecto a sus participaciones en el ingreso y población total. El índice de Theil para estos datos agrupados (T^a) se define como

$$(52) \quad T^a = \sum_g y_g \text{Log} \frac{y_g}{n_g}$$

Por lo tanto constatamos que

$$(53) \quad T^a = B_g$$

Si se contara con la información a un nivel individual se podría computar el verdadero índice de Theil (T) para ese conjunto de datos

$$(54) \quad T = \sum_{u=1}^n y_u \text{Log} \frac{y_u}{1/N} = \sum_g y_g \text{Log} \frac{y_g}{n_g} + \sum_g y_g \sum_{u \in S_g} \frac{y_{gu}}{y_g} \text{Log} \frac{\frac{y_{gu}}{y_g}}{1/N_g}$$

$$T = B_g + W_g$$

Por lo tanto

$$(55) \quad T - T^a = W_g = \sum_g y_g T_g \geq 0$$

Lo que implica que el índice de Theil agrupado será siempre menor o igual que el desagregado. Esto se debe a que el índice agrupado, al no contar con información a nivel individual dentro de cada grupo, supone que la dispersión del ingreso dentro de ellos es nula.

b) Datos agrupados y análisis de descomposición

Una forma en que se presenta el problema de los datos agrupados es la siguiente. Existe una variable estratificadora j que puede tomar \bar{j} valores distintos generando \bar{j} grupos S_j con N_j individuos cada uno. Sin embargo, no se cuenta con las participaciones en el ingreso de cada uno de los individuos de los distintos grupos, sino que sólo se conoce para cada grupo S_j la distribución de sus miembros, entre los distintos tramos de ingresos (n_{jt}) y la participación en el ingreso total de cada uno de los tramos (y_{jt}) (es decir, el ingreso promedio en cada tramo) en que el subíndice t se refiere a los distintos tramos de ingresos en cada grupo S_j . En estas circunstancias

$$(56) \quad T^a = \sum_j \sum_t y_{jt} \text{Log} \frac{y_{jt}}{n_{jt}} = \sum_j y_j \text{Log} \frac{y_j}{n_j} + \sum_j y_j \sum_t \frac{y_{jt}}{y_j} \text{Log} \frac{y_{jt} y_j}{n_{jt} n_j}$$

$$(57) \quad T^a = B_j + W_j^a$$

$$(58) \quad W_j^a = \sum_j y_j T_j^a$$

En tanto que

$$(59) \quad T = \sum_j \sum_t \sum_u y_{jtu} \text{Log} \frac{y_{jtu}}{1/N} = \sum_j y_j \text{Log} \frac{y_j}{n_j} + \sum_j y_j \sum_t \frac{y_{jt}}{y_j} \text{Log} \frac{y_{jt} y_j}{n_{jt} n_j} +$$

$$\sum_j \sum_t y_{jt} \sum_u \frac{y_{jtu}}{y_{jt}} \text{Log} \frac{y_{jtu}}{1/N_{jt}}$$

$$(60) \quad T = B_j + W_j^a + W_{jt}$$

$$(61) \quad W_j = W_j^a + W_{jt}$$

Por lo tanto observamos que

$$(62) \quad T - T^a = W_{jt} = \sum_j \sum_t y_{jt} T_{jt} \geq 0$$

Esto indica que el índice agrupado subestima al verdadero índice de Theil de esta población ya que el primero al tomar los datos agrupados prescinde de parte de la desigualdad en la distribución de los ingresos: la desigualdad dentro de los distintos tramos de ingresos de cada uno de los grupos S_j , mientras mayor sea la dispersión del ingreso dentro de cada uno de estos grupos mayor sea la discrepancia entre el índice agrupado y el verdadero índice.

Sin embargo, la contribución absoluta a la desigualdad total de la variable j (B_j) no se ve afectada por el uso de datos agrupados siendo idéntica a la que se hubiera obtenido usando datos no agrupados. Esto no ocurre con su contribución como proporción de la desigualdad total.

$$(63) \quad \frac{B_j}{T^a} - \frac{B_j}{T} = \frac{B_j}{T} \frac{T}{T^a} - \frac{B_j}{T} = \frac{B_j}{T} \left[\frac{T}{T^a} - 1 \right]$$

$$(64) \quad \frac{\frac{B_j}{T^a} - \frac{B_j}{T}}{\frac{B_j}{T}} = \frac{T - T^a}{T^a} = \frac{W_{jt}}{T^a} \gg 0$$

Por lo tanto, al trabajar con datos agrupados, la contribución porcentual de la variable j a la desigualdad total excede la verdadera contribución

en un $\frac{W_{jt}}{T^a}$ por ciento. Mientras mayor sea la dispersión del ingreso al interior de cada uno de los tramos de ingresos con respecto a la dispersión entre tramos de ingresos, mayor será el sesgo porcentual.

Recordando la expresión (62) vemos que estimando bajo "supuestos razonables" el nivel de desigualdad dentro de los distintos tramos de ingresos se puede estimar W_{jt} y obtener así una estimación del sesgo porcentual identificado en (64).

En el caso general en que se clasifica la población de acuerdo a R características se tiene que

$$(65) \quad T^a = \sum_{C_1} \sum_{C_2} \dots \sum_{C_R} \sum_t y_{C_1, C_2 \dots C_R, t} \text{Log} \frac{y_{C_1, C_2 \dots C_R, t}}{n_{C_1, C_2, C_R, t}} =$$

$$\sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2 \dots C_R} \text{Log} \frac{y_{C_1, C_2 \dots C_R}}{n_{C_1, C_2 \dots C_R}} +$$

$$\sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2 \dots C_R} \sum_t \frac{y_{C_1, C_2 \dots C_R, t}}{y_{C_1, C_2 \dots C_R}} \text{Log} \frac{\frac{y_{C_1, C_2 \dots C_R, t}}{y_{C_1, C_2 \dots C_R}}}{\frac{n_{C_1, C_2 \dots C_R, t}}{n_{C_1, C_2 \dots C_R}}}$$

$$(66) \quad T^a = B_{1, 2 \dots R} + W_{1, 2 \dots R}^a$$

$$(67) \quad W_{1, 2 \dots R}^a = \sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2 \dots C_R} T_{1, 2 \dots R}^a$$

En tanto que

$$(68) \quad T = \sum_{C_1} \sum_{C_2} \dots \sum_{C_R} \sum_t \sum_u y_{C_1, C_2 \dots C_R, t, u} \text{Log} \frac{y_{C_1, C_2 \dots C_R, t, u}}{\frac{1}{N}} =$$

$$\sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2 \dots C_R} \text{Log} \frac{y_{C_1, C_2 \dots C_R}}{n_{C_1, C_2 \dots C_R}} +$$

$$\sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2 \dots C_R} \sum_t \frac{y_{C_1, C_2 \dots C_R, t}}{y_{C_1, C_2 \dots C_R}} \text{Log} \frac{\frac{y_{C_1, C_2 \dots C_R, t}}{y_{C_1, C_2 \dots C_R}}}{\frac{n_{C_1, C_2, C_R, t}}{n_{C_1, C_2, C_R}}}$$

$$\sum_{C_1} \sum_{C_2} \dots \sum_{C_R} \sum_t y_{C_1, C_2 \dots C_R, t} \sum_u \frac{y_{C_1, C_2 \dots C_R, t, u}}{y_{C_1, C_2 \dots C_R, t}} \text{Log} \frac{\frac{y_{C_1, C_2 \dots C_R, t, u}}{y_{C_1, C_2 \dots C_R, t}}}{\frac{1}{N} \frac{n_{C_1, C_2 \dots C_R, t}}{n_{C_1, C_2 \dots C_R}}}$$

$$(69) \quad T = B_{1,2..R} + W_{1,2..R}^a + W_{1,2..R,t}$$

Por lo tanto

$$(70) \quad T - T^a = W_{1,2..R,t} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_R} y_{C_1, C_2 \dots C_R} T_{1,2..R,t} \geq 0$$

$$(71) \quad \frac{\frac{B_{1,2..R}}{T^a} - \frac{B_{1,2..R}}{T}}{\frac{B_{1,2..R}}{T}} = \frac{T - T^a}{T^a} = \frac{W_{1,2..R,t}}{T^a} \geq 0$$

Vemos que las contribuciones individuales brutas marginales y conjuntas medidas como proporción de la desigualdad total incurren en un sesgo

de $\frac{W_{1,2..R,t}}{T^a}$ por ciento, al computarse con datos agrupados. Sin

embargo, si se expresan las distintas contribuciones de las R variables (individuales brutas marginales y conjuntas de orden menor a R) como porcentajes de la contribución conjunta de todas ellas, estas contribuciones relativas no se verán afectadas por el uso de datos agrupados, facilitando por tanto la comparación de resultados obtenidos de datos agrupados y no agrupados. Esto debido a que el uso de datos agrupados, introduce un sesgo en las contribuciones medidas como porcentaje de la desigualdad total pero no afecta los valores absolutos de estas contribuciones y por tanto tampoco los valores relativos entre sí.

Con respecto a las interacciones. Dado que las interacciones son combinaciones lineales de contribuciones conjuntas de distinto orden, las cuales no se ven afectadas por el uso de datos agrupados, el valor absoluto de las interacciones tampoco se verá afectado por el uso de datos agrupados.

En síntesis, gran parte del análisis puede ser llevado a cabo con datos agrupados en forma tal que los resultados así obtenidos sean comparable con los resultados obtenidos de datos no agrupados. Sólo en los casos en que los resultados de contribuciones o interacciones se expresen como porcentaje de la desigualdad total se hace necesario corregir el sesgo introducido por el uso de datos agrupados e identificado en la expresión (71).

VI. Ordenación de las variables de acuerdo con la contribución a la explicación de la desigualdad de los ingresos

Del método de descomposición de Theil, no se deriva un criterio único que permita ordenar las variables de acuerdo a la importancia de su contribución a la explicación de la desigualdad de los ingresos.

En esta sección se definirán y analizarán brevemente cuatro criterios alternativos tendientes a esa ordenación y que obedecen a propósitos distintos.

- a) Un primer criterio para ordenar las variables es hacerlo de acuerdo a sus contribuciones individuales brutas. Este criterio ordena las variables de acuerdo a la contribución que tienen al ser tomadas en forma aisladas y puede por lo tanto incluir el impacto de otras variables ejercido a través de su asociación estadística con la variable en cuestión. Este criterio tiene la ventaja de que la contribución individual bruta de las variables es independiente del universo de variables consideradas y por tanto, la ordenación resultante no se verá afectada por la inclusión de una variable adicional a este universo.
- b) Un segundo criterio consiste en ordenar las variables de acuerdo con la secuencia con que maximizan la contribución conjunta de grupos con un número creciente de variables en la primera variable seleccionada sería aquella con la mayor contribución individual bruta. La segunda variable, aquella cuya contribución marginal dada la ya seleccionada sea máxima, es decir, aquella que maximice la contribución conjunta del grupo de dos variables que se generará al seleccionar la segunda. La tercera variable seleccionada sería aquella con la mayor contribución marginal dadas las dos primeras y así sucesivamente ir seleccionando las variables de acuerdo con sus contribuciones marginales dadas las ya seleccionadas. Este criterio considera parcialmente las interacciones entre variables puesto que para seleccionar la segunda y demás variables considera la interacción entre las candidatas a ocupar estos lugares y las ya seleccionadas. Esto hace que la ordenación de las variables de acuerdo a este criterio dependa del universo de variables consideradas. La inclusión de una nueva variable

a este universo puede alterar la ordenación de las ya existentes obtenida previamente. Este criterio, al ordenar las variables en forma tal de maximizar la contribución conjunta de conjuntos con un número creciente de variables, debería orientar la decisión respecto a cuáles variables y en qué orden debieran ser incluidas en el análisis cuando existe una restricción respecto al número máximo de variables a incluir.

c) Un tercer criterio para ordenar las variables es secuencial como el anterior, pero en sentido inverso. Consiste en seleccionar como primera variable aquella que más reduciría el poder explicativo del conjunto inicial de variables si fuera retirada. Es decir, aquella cuya contribución marginal dadas las demás sea máxima. Seleccionar como segunda variable aquella que, dada el retiro de la primera, más reduciría el poder explicativo del conjunto de variables restantes. Es decir, ir seleccionando las variables de acuerdo con su contribución marginal dado el conjunto de variables reducido por el retiro de las anteriores. Este criterio también considera parcialmente las interacciones entre variables ya que para ordenarlas considera las interacciones entre las candidatas a un determinado lugar y las que no han sido previamente seleccionadas y retiradas del análisis. Esto hace que la ordenación de las variables de acuerdo a este criterio también dependa del universo de variables consideradas. Esta ordenación se puede ver afectada por la inclusión de una nueva variable al universo. La ordenación de las variables en sentido opuesto al propuesto por este criterio tiende a minimizar la pérdida de poder explicativo que se origina al retirar variables del análisis y por tanto, este criterio debería orientar la decisión respecto de cuáles variables y en qué orden se debería prescindir si esto fuera necesario.

d) Finalmente un último criterio para ordenar las variables consiste en hacerlo de acuerdo con sus contribuciones marginales dadas todas las demás. Este criterio considera íntegramente las interacciones entre las variables y por tanto la ordenación resultante de él también depende del universo de variables consideradas. Este criterio, al basarse en la contribución marginal, tiene la ventaja de que la

contribución marginal de una variable se computa controlando el efecto de todas las demás, y por tanto representa el incremento en el grado de explicación exclusivamente atribuible a esa variable.

Es importante destacar que en la medida que un criterio considere parcial o totalmente las interacciones entre las variables, la ordenación resultante dependerá del universo de variables consideradas. Sin embargo, mientras mayor sea la proporción de la desigualdad total explicada por el universo considerado, menor la probabilidad de que la inclusión de una variable adicional altere la ordenación.

A partir de las expresiones (48), (49) y (51) se desprende que si todas las interacciones entre las variables fueran nulas, la contribución individual de una variable i sería idéntica a cualquier contribución marginal de esta variable cualesquiera sea el grupo de variables controladas. Por lo tanto, si todas las interacciones fueran nulas, los cuatro criterios de ordenamiento propuestos anteriormente coincidirían perfectamente. Sin embargo, en la medida que existan interacciones, las ordenaciones de las variables resultantes de estos criterios pueden diferir diametralmente.

La determinación de la influencia relativa de las variables sobre la distribución agregada del ingreso puede, sin embargo, ocultar la existencia de sistemas de influencia e interacciones ampliamente diferentes para distintos segmentos de la población total. En tales casos, resulta conveniente reproducir el análisis de descomposición de la desigualdad en el interior de cada uno de esos segmentos. Por ejemplo, la descomposición de las desigualdades urbanas y de las rurales, o la descomposición de las desigualdades en el interior de cada categoría del empleo, pueden revelar con más profundidad los mecanismos que determinan las desigualdades del ingreso a nivel total.

Apéndice Matemático

1. Límite inferior y superior del índice de Theil

i) Cota inferior

Dado que $\text{Log } x \leq x-1$ se tiene que

$$-\sum_i y_i \text{Log} \frac{y_i}{n_i} = \sum_i y_i \text{Log} \frac{n_i}{y_i} \leq \sum_i y_i \left(\frac{n_i}{y_i} - 1 \right) = \sum_i n_i - \sum_i y_i = 1 - 1 = 0$$

Por lo tanto

$$\sum_i y_i \text{Log} \frac{y_i}{n_i} \geq 0$$

ii) Cota superior

$$\sum_i y_i \text{Log} \frac{y_i}{1/N} = \sum_i \text{Log} \left(y_i^{y_i} \right) + \text{Log } N$$

pero

$$y_i \leq 1 \Rightarrow y_i^{y_i} \leq 1 \Rightarrow \text{Log } y_i^{y_i} \leq 0$$

Por lo tanto

$$\sum_i y_i \text{Log} \frac{y_i}{1/N} \leq \text{Log } N$$

2. Para probar que la expresión (24) es no-negativa, se usa la misma desigualdad mencionada anteriormente

$$-\sum_i \sum_j y_{ij} \text{Log} \frac{y_{ij}}{y_i y_j} = \sum_i \sum_j y_{ij} \text{Log} \frac{y_i y_j}{y_{ij}} \leq \sum_i \sum_j y_{ij} \left(\frac{y_i y_j}{y_{ij}} - 1 \right) =$$

$$\sum_i \sum_j y_i y_j - \sum_i \sum_j y_{ij} = 1 - 1 = 0$$

Por lo tanto

$$\sum_i \sum_j y_{ij} \text{Log} \frac{y_{ij}}{y_i y_j} \geq 0$$

3. Demostración de la descomposición del índice de Theil presentado en la expresión (32).

La demostración se hará aplicando el método de inducción sobre el número de variables consideradas: R

Si R = 1

$$\begin{aligned}
 (1) \quad T &= \sum_{u=1}^n y_u \operatorname{Log} \frac{y_u}{1/N} = \sum_{C_1=1}^{\bar{C}_1} \sum_{u \in C_1} y_{C_1,u} \operatorname{Log} \frac{y_{C_1,u}}{1/N} = \\
 &= \sum_{C_1} \sum_u y_{C_1,u} \operatorname{Log} \frac{y_{C_1}}{n_{C_1}} \frac{y_{C_1,u}/y_{C_1}}{1/N/n_{C_1}} = \sum_{C_1} y_{C_1} \operatorname{Log} \frac{y_{C_1}}{n_{C_1}} + \\
 &= \sum_{C_1} y_{C_1} \sum_{u \in C_1} \frac{y_{C_1,u}}{y_{C_1}} \operatorname{Log} \frac{y_{C_1}}{1/N_{C_1}} = B_1 + W_1
 \end{aligned}$$

Por lo tanto, la descomposición es válida para una variable. Supongamos que la descomposición es válida para (q-1) variables entonces

$$(2) \quad T = B_1 + B_2^1 + \dots + B_{q-1}^{1,2,\dots,(q-2)} + W_{1,2,\dots,(q-1)}$$

Peró el residuo $W_{1,2,\dots,(q-1)}$ se definía en la expresión (36) como

$$\begin{aligned}
 (3) \quad W_{1,2,\dots,(q-1)} &= \sum_{C_1} \sum_{C_2} \dots \sum_{C_{(q-1)}} y_{C_1, C_2, \dots, C_{q-1}} \\
 &= \sum_{u \in S_{C_1, C_2, \dots, C_{(q-1)}}} \frac{y_{C_1, C_2, \dots, C_{(q-1)}, u}}{y_{C_1, C_2, \dots, C_{(q-1)}}} \operatorname{Log} \frac{y_{C_1, C_2, \dots, C_{(q-1)}, u} y_{C_1, C_2, \dots, C_{(q-1)}}}{1/N_{C_1, C_2, \dots, C_{(q-1)}}}
 \end{aligned}$$

Esto puede descomponerse en las diversas sumas correspondientes a los grupos más pequeños $S_{C_1, C_2, \dots, C_{(q-1)}} C_q$

$$(4) \quad W_{1,2..(q-1)} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_{(q-1)}} \sum_{C_q} \sum_{u \in S_{C_1, C_2 \dots C_q}} y_{C_1, C_2 \dots C_{(q-1)}, C_q, u}$$

$$\text{Log} \frac{y_{C_1, C_2 \dots C_{(q-1)}, C_q, u}}{1 / N_{C_1, C_2 \dots C_{(q-1)}}} =$$

$$\sum_{C_1} \sum_{C_2} \dots \sum_{C_q} \sum_{u \in S_{C_1, C_2 \dots C_q}} y_{C_1, C_2 \dots C_{(q-1)}, C_q, u}$$

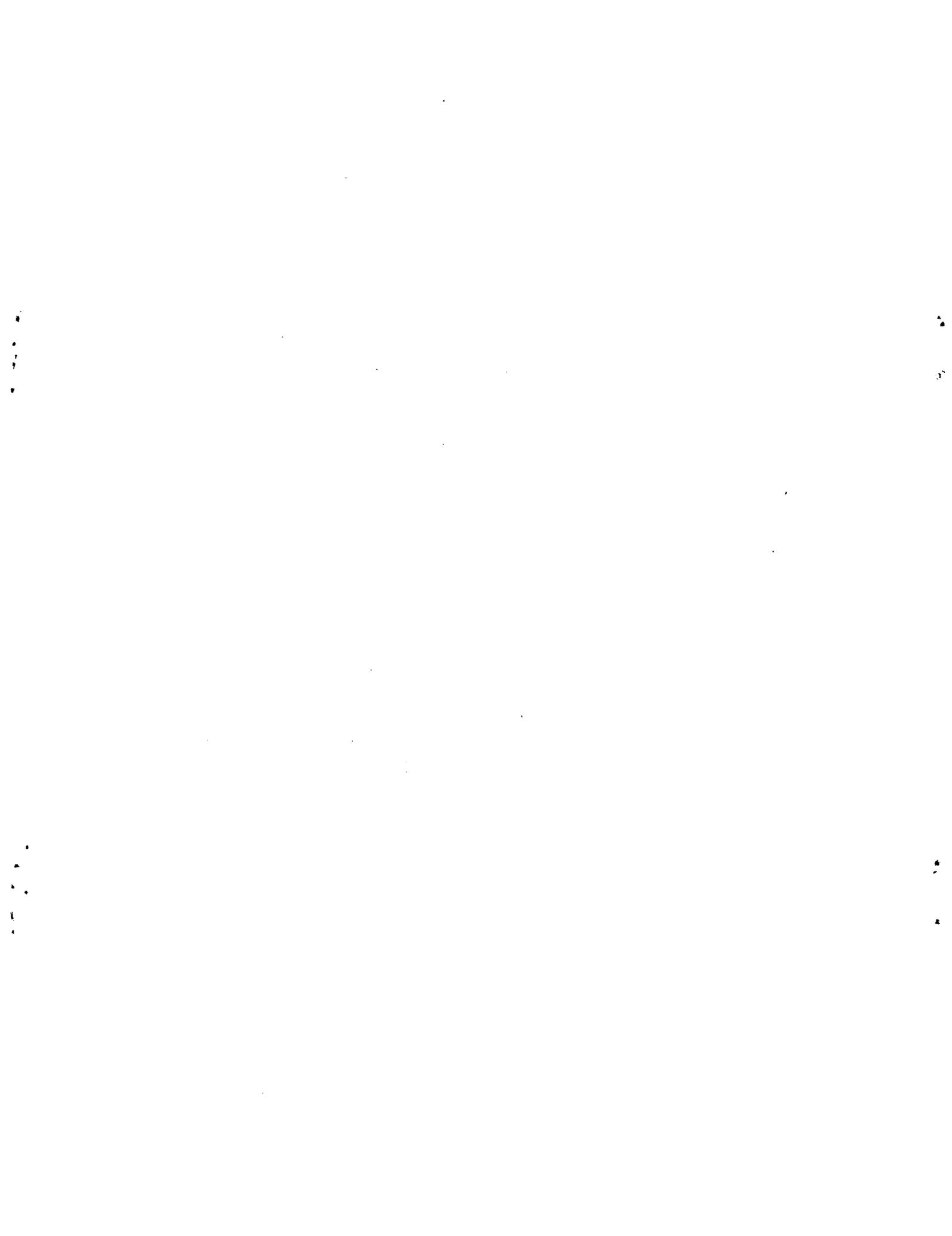
$$\text{Log} \frac{\frac{y_{C_1, C_2 \dots C_q}}{y_{C_1, C_2 \dots C_{(q-1)}}}}{\frac{n_{C_1, C_2 \dots C_q}}{n_{C_1, C_2 \dots C_{q-1}}}} = \sum_{C_1} \sum_{C_2} \dots \sum_{C_{q-1}} y_{C_1, C_2 \dots C_{(q-1)}}$$

$$\text{Log} \frac{\frac{y_{C_1, C_2 \dots C_q}}{y_{C_1, C_2 \dots C_{q-1}}}}{\frac{n_{C_1, C_2 \dots C_q}}{n_{C_1, C_2 \dots C_{q-1}}}} + \sum_{C_1} \sum_{C_2} \dots \sum_{C_q} \sum_{u \in S_{C_1 \dots C_q}} y_{C_1, C_2 \dots C_q, u}$$

$$\text{Log} \frac{y_{C_1, C_2 \dots C_q, u}}{1 / N_{C_1, C_2 \dots C_q}} = B_q^{1,2..(q-1)} + W_{1,2..q}$$

Por lo tanto la expresión (3) puede escribirse como

$$(5) \quad T = B_1 + B_2^1 + B_3^{1,2} + \dots + B_q^{1,2..(q-2)} + B_q^{1,2..(q-1)} + W_{1,2..q}$$



REFERENCIAS BIBLIOGRAFICAS

- Chiswick, C. (1976) Aplication of the Theil Index To Income
Inequality Working Paper Series B-2, Dev.
Research Center, World Bank.
- Fishlow, A (1972) "Brazilian Size Distribution of Income"
American Economic Review, Mayo 1972.
- Herfindall, O.C. (1950) Concentration in Steel Industry. Unpublished
Ph.D. Dissertation, Columbia University.
- Hirschman, A.O. (1945) National Power and The Structure of Foreign
Trade University of California Press.
- Rice, S.R. (1928) Quantitative Methods in Politics, New York,
A.A. Knopf Inc.
- Shannon, C.E. (1948) "A Mathematical Theory of Communication"
Bell System Technical Journal, Vol. 27.
- Theil, H. (1967) Economics and Information Theory
North Holland
- (1972) Statistical Decomposition Analysis
North Holland
- van Ginneken, W. (1975) "Análisis de descomposición del Índice de
Theil aplicado a la distribución del ingreso
familiar en México", Demografía y Economía,
Nº 25, Vol IX, Nº 1, 1975.

