

Cátedra: Sautu

Materia: Método I

Carrera:

Autor: ANÁlisis Bivariado
GARCÍA Fernando Ap. 7

- 53 por 100, con una desviación típica del 3 por 100. ¿Hasta qué punto se puede afirmar que el grado de apoyo de los electores al candidato en la campaña actual es diferente del manifestado en el pasado, al nivel de significación de a) 0,05 y b) 0,005?
6. De una muestra de 500 votantes elegidos aleatoriamente en una comunidad, el 55 por 100 de ellos son partidarios de un determinado candidato. Hallar los límites de confianza del a) 95 por 100 y b) 99 por 100 para la proporción de todos los votantes que son partidarios de dicho candidato.
 7. ¿Qué tamaño de muestra debería tomarse en el caso del ejercicio anterior para que la confianza de que el candidato salga elegido (es decir, obtenga el 50 por 100 o más de votos) fuera del 95 por 100? y del 99 por 100?
 8. En un estudio sobre la ideología de los trabajadores de una gran empresa, se pidió a 100 trabajadores elegidos al azar que se posicionasen en una escala de preferencia política que tiene un recorrido del 1 al 10 (1=extrema izquierda; 10=extrema derecha). El valor medio de los posicionamientos de los trabajadores fue $\bar{X}=4,2$, con una desviación típica de 0,04. Estimar el valor medio del posicionamiento ideológico de todos los trabajadores, con un intervalo de confianza de a) 95 por 100 y b) 99 por 100.

BIBLIOGRAFIA

- BLALOCK, Hubert M.: *Social Statistics*, New York, McGraw-Hill, 1960.
 COCHRAN, W. G.: «Some methods for strengthening the common χ^2 test», *Biometrics*, 10, 1954, pág. 417-451.
 DOMÉNECH i MASSONS, J. M.: *Bioestadística. Métodos Estadísticos para Investigadores*, Barcelona, Herder, 1977.
 LABOVITZ, Sanford: «The Nonutility of Significance Tests: The Significance of Test of Significance Reconsidered», *Pacific Sociological Research*, vol. 13, núm. 3, 1970, págs. 141-147.
 LOETHER, N. J., y D. G. McTAVISH: *Inferential Statistics for Sociologists*, Boston, Allyn & Bacon, 1974.
 MORRISON, D. E., y R. E. HENKEL: *The Significance Test Controversy*, Chicago, Aldine, 1970.
 SÁNCHEZ-CRESPO, J. L.: *Principios elementales del muestreo y estimación de proporciones*, Madrid, INE, 1971.
 SIEGEL, Sidney: *Nonparametric Statistics for the Behavioral sciences*, New York, McGraw-Hill, 1956.

Capítulo 7

ESTADÍSTICA DESCRIPTIVA BIVARIABLE: CARACTERÍSTICAS DE UNA ASOCIACIÓN BIVARIABLE

Nuestro objetivo en el presente capítulo es el estudio de las características de las distribuciones bivariadas o clasificaciones cruzadas de dos variables. Con ello adelantamos un nuevo paso en nuestro recorrido por el camino que nos va mostrando, en su creciente complejidad, la diversidad de las técnicas estadísticas utilizadas en la investigación sociológica. En la estadística descriptiva univariada comenzamos con una distribución de frecuencias y, a partir de ella, estudiamos una serie de medidas resumen que nos permitieron lograr números índices, de gran utilidad para la descripción de los datos sociológicos. Al mismo tiempo, se desarrollaron una serie de medidas para cada uno de los rasgos definitorios de una distribución; esto es, la tendencia central, la forma y la variabilidad o dispersión. En el presente capítulo, nuestro objetivo es similar, pero, si cabe, más interesante, pues nos vamos acercando más al tipo de tarea que con mayor frecuencia se realiza en la realidad de la investigación sociológica, esto es, el estudio de las condiciones que influyen en la distribución de una variable.

7.1. DISTRIBUCIONES BIVARIABLES: UN EJEMPLO

Los sociólogos que estudian las actitudes políticas de la población vienen utilizando, desde hace años, en las encuestas de opinión una escala de preferencia ideológica izquierda-derecha. En la entrevista se pide al entrevistado que se sitúe en una de las casillas que componen una escala que, según el tipo de estudio, va del 1 al 7 o del 1 al 10, correspondiendo el extremo 1 a la preferencia por la extrema izquierda y el extremo 7, ó 10, a la preferencia por la extrema derecha. Como señalan los autores del estudio *La Conciencia Regional en España* (J. Jiménez Blanco *et al.*, 1977, pág. 88), «se trata de un planteamiento de extrema simplicidad, en el que se traspone la dimensión ideológica del espectro político a una dimensión geométrica en el plano». Pues bien, los porcentajes de autoubicación obtenidos de la población española en el referido estudio fueron los siguientes.

TABLA 7.2

Escala Izquierdo-derecha entre la población clasificada según nivel de religiosidad

Notese que estímulos ligeramente diferentes a estos conclusiones traves de la compa-
nía de porcentajes, en lugar de frecuencias absolutas, ya que el nü-
mero de miembros que corresponde a cada uno de los cinco grupos
considerados es diferente entre si. Por eso, el uso de porcentajes es ta-
cticamente más efectivo para realizar comparaciones válidas, tal como se estudió en
el capítulo 2.

Lo que hemos hecho hasta ahora ha sido comparar la distribución
de los valores de la escala izquierda-derecha entre diversas distribucio-
nes universitables. La variable dependiente es común a cada una de las
distribuciones, y cada tabla se diferencia porque agrupa a los individuos
que pertenecen a la misma categoría. Pero una forma más
eficaz y rápida de obtener conclusiones válidas bajo estas condiciones
sería la de combinar las cinco tablas separadas en una sola tabla, tal
como se hace a continuación.

La interpretación de estos datos pude ser muy diversa, según el número de casillas que se asignan a cada postura ideológica. Haciendo una «lectura desdoblada» de las casillas que el centro las ensayan los autores, como ensayarán las primeras cuatro a la derecha, quedando para el centro las casillas cinco y seis. Ahora bien, los autores quedan para la izquierda y las últimas cuatro a la derecha, quedando para el centro casillas a la izquierda de las últimas cuatro. La diferencia es que la izquierda no deseaba tan solo conocer la distribución global de la población en función de otras variables relevantes. Así, a partir de la consideración en relación a lo largo de dicha escala, sino que, además, indagaron sobre la distribución de la población en función de la religiosidad. En la tabla 7.2 aparece el conjunto de tales distribuciones. Se trata de cinco distribuciones univariadas de datos de observación que se producen en la distribución de la población entre las diversas posturas ideológicas, para los distintos grupos considerados. Una parte imágenes, que se declaran más diferentes en materia religiosa, o que portante de los que se declaran más semejantes no prácticas religiosas, una parte imágenes católicas que todo lo contrario ocurre entre los que se consideran a sí mismos católicos no prácticas religiosas, se autodefinen en la derecha, que tienden a situarse en el centro o, sobre todo, muy buenas católicas, mientras que todo lo contrario ocurre entre los que se consideran a sí mismos católicas no prácticas religiosas, se autodefinen en la derecha. Los católicos prácticas o no muy prácticas se sitúan preferentemente en las casillas 5 y 6, correspondientes al centro.

TABLA 7.1

Izquierda.	Escala	%	/	
1	127	2	3	3
2	190	3	3	3
3	381	4	4	4
4	444	5	5	5
5	1.522	6	6	6
6	888	7	7	7
7	381	8	8	8
8	444	9	9	9
9	254	14	14	14
10	317	24	24	24
Derecha.	NS/NC	21	21	21
	1.334	5	5	5
	6.342	100	100	100

Porcentajes de autoubicación en un espacio político abstracto

TABLA 7.3

Distribución porcentual de la escala izquierda-derecha según el nivel de religiosidad de la población española

Escala izquierda-derecha	Total	Muy buenos católicos	Católicos practicantes	Católicos no muy practicantes	Católicos no practicantes	Indiferentes
Izqda. (1-4) ...	18	5	8	16	34	52
Centro (5-6) ...	38	30	43	35	37	28
Dcha. (7-10) ...	22	39	26	18	11	6
NS/NC ...	21	26	23	21	18	14
Total ...	100	100	100	100	100	100
N =	(6.342)	(903)	(2.366)	(1.466)	(859)	(677)

FUENTE: Ver tablas 7.1 y 7.2

Esta tabla nos permite comparar los diversos grupos entre sí, y cada uno de ellos con la media nacional, de una forma más rápida y eficaz, a la vez que ofrece un excelente resumen de la información que se contiene en las dos tablas anteriores, en forma del tipo de relación que se produce entre dos variables. Es el tipo de tabla que se conoce con el nombre de *distribución porcentual bivariante*, porque permite examinar la distribución porcentual de una variable (la variable dependiente) dentro de las diferentes categorías de otra variable (la variable independiente). Precisamente las ideas que subyacen debajo de tales clasificaciones cruzadas constituyen las bases del análisis empírico en la sociología, ya que es a través de dicho análisis como se trata de formular y contrastar el tipo de relación existente entre las variables, así como las condiciones en que se produce.

Una distribución bivariante, tal como la que se presenta en la tabla 7.3, permite no sólo el examen directo de la distribución global de una variable dependiente, sino también las condiciones que se supone influyen en la manera en que se distribuye dicha variable. Por lo que se refiere al caso de la ideología política, la teoría sugiere que, bajo ciertas condiciones, las posturas ideológicas de izquierda, centro y derecha se distribuirán de forma diferente que bajo otras condiciones. En el caso concreto de los datos que se incluyen en la tabla 7.3, tales condiciones corresponden a los diferentes niveles de religiosidad, aunque, como resulta obvio, el sociólogo puede pensar en otras condiciones que también pueden influir en la preferencia ideológica, tales como la edad, el sexo, la clase social, el lugar de residencia, etc. Como se ha dicho otras veces, se puede imaginar que una de las tareas de toda disciplina científica es la de buscar los tipos de condiciones que contribuyen a

mejor predecir y explicar el nivel de algún tipo de fenómeno (Loether y McTavish, 1974, pág. 174). Una tabla bivariante como la anterior, en realidad, pone en relación una serie de distribuciones condicionales con una distribución global de una variable dependiente.

En lo que resta de capítulo vamos a estudiar las principales características de las distribuciones bivariadas o clasificaciones cruzadas de dos variables, y comenzaremos dicho estudio estableciendo algunas reglas prácticas sobre la forma de presentar y leer correctamente, desde un punto de vista tanto teórico como metodológico, una tabla bivariante.

7.2. PRESENTACIÓN Y ANÁLISIS DE UNA TABLA BIVARIANTE

La tabulación cruzada y comparativa de dos variables da lugar a una tabla compuesta de filas y columnas, utilizando las categorías de cada variable para designar, respectivamente, las filas y las columnas. Se suele seguir la convención de situar la variable dependiente, cuando la hubiere, en las filas, y la variable independiente en las columnas.

Con el fin de ilustrar la forma en que se construye una tabla, supongamos que disponemos de 15 puntuaciones correspondientes a los

TABLA 7.4

Ejemplo de una tabla de frecuencias bivariante

Puntuaciones	Sexo	Escala izquierda-derecha	Puntuaciones (cont.)	Sexo	Escala izquierda-derecha
1	v	I	8	m	C
2	m	D	9	m	D
3	v	C	10	v	D
4	v	C	11	v	I
5	m	D	12	m	I
6	v	I	13	v	C
7	m	D	14	m	C
			15	m	I

Tabla de preferencia ideológica por sexo

	Escala izqda.-dcha	Sexo		Total filas
		v	m	
v — varón				
m — mujer				
I — izquierda	I	3	2	5
C — centro	C	3	2	5
D — derecha	D	1	4	5
Total columnas		7	8	15

Las comparaciones que se desean comparar. En tal caso surge la pregunta: ¿que comparación que se deseará tener en las diferentes filas?

Existen una regla sencilla, universalmente aceptada, que se utiliza como guía para responder a la anterior pregunta. Dicha regla plantea que no se tratará en cuanta, al aplicarla de «causa y efecto», que no se tratará en la variable independiente que se considera del otro factor, sino de que factor causará o de la variable independiente en el sentido del factor «causal» o de la variable independiente.

Comúnmente se considera la causalidad de resolver un problema de causalidad real, del segundo factor. Por esa razón, algunas autoras prefieren hablar de variables independientes y dependientes, en cuyo caso se dice que los factores causales son las variables que se consideran en la variable independiente.

Variables independientes y dependientes, en cuyo caso se dice que los factores causales son las variables que se consideran en la variable independiente.

Establecida una tabla 2×3 (dos columnas \times tres filas), se forman seis celulas en las que se escriben las correspondientes frecuencias que aparecen en la distribución global de puntajes y de los datos. De esta forma se calculan dos subtotalares (el subtotal de las filas y el subtotal de las columnas) y un total global.

Para una tabla 2×2 , se pueden simbolizar las operaciones que hemos realizado en el ejemplo anterior del siguiente modo:

Habituadamente, las tablas bivariadas se presentan con números que reflejan porcentajes en lugar de frecuencias absolutas. Con ello se facilita la realización de comparaciones numéricas entre las distribuciones de las tablas bivariadas que tienen las mismas características que las tablas bivariadas que se presentan en una tabla.

7.2.a. Cálculo de los porcentajes en una tabla

El valor de la frecuencia correspondiente a cada celdilla se simboliza mediante n_{ij} , en donde el subíndice i indica el número de la fila y el subíndice j indica el número de la columna. Esto es, que la frecuencia de cada celdilla indica el número de casos que cumplen las condiciones simultáneamente de que la fila i y la columna j se observan.

	Columna 1	Columna 2	Total filas	
Fila 1	n_{11}	n_{12}	n_1	Total columnas
Fila 2	n_{21}	n_{22}	n_2	
				n

Establecida una tabla 2×3 (dos columnas \times tres filas), se forman seis celulas en las que se escriben las correspondientes frecuencias que aparecen en la distribución global de puntajes y de los datos. De esta forma se calculan dos subtotalares (el subtotal de las filas y el subtotal de las columnas) y un total global.

Para una tabla 2×2 , se pueden simbolizar las operaciones que hemos realizado en el ejemplo anterior del siguiente modo:

El valor de la frecuencia correspondiente a cada celdilla se simboliza mediante n_{ij} , en donde el subíndice i indica el número de la fila y el subíndice j indica el número de la columna. Esto es, que la frecuencia de cada celdilla consta de los subtotalares que cumplen las condiciones simultáneamente de que la fila i y la columna j se observan.

TABLA 7.5

Ilustración de las diferentes formas en que se pueden calcular los porcentajes en una tabla bivariante

Distribución de las frecuencias de la tabla 7.4.

Escala izqda.-dcha.	Sexo		
	Varón	Mujer	Total
I	3	2	5
C	3	2	5
D	1	4	5
Total	7	8	15

7.5.a. Cálculo de porcentajes tomando como base los totales de las columnas.

Escala izqda.-dcha.	Sexo		
	Varón %	Mujer %	Total %
I	42,9	25,0	33,3
C	42,9	25,0	33,3
D	14,2	50,0	33,3
Total	100	100	100

7.5.b. Cálculo de porcentajes tomando como base los totales de las filas.

Escala izqda.-dcha.	Sexo		
	Varón %	Mujer %	Total %
I	60,0	40,0	100
C	60,0	40,0	100
D	20,0	80,0	100
Total	46,7	53,3	100

7.5.c. Cálculo de porcentajes tomando como base la frecuencia total N.

Escala izqda.-dcha.	Sexo		
	Varón %	Mujer %	Total %
I	20,0	13,3	33,3
C	20,0	13,3	33,3
D	6,7	26,7	33,3
Total	46,7	53,3	100

porcentajes de «izquierda» entre hombres y mujeres es 17,9 por 100 ($42,9 - 25,0 = 17,9$ por 100). Este valor se llama *epsilón*, y se simboliza mediante la letra griega ϵ (Loether y McTavish, *op. cit.*, pág. 178). En tablas con más de dos columnas se pueden calcular porcentajes de contraste o *epsilones* para cada par de columnas y entre las correspondientes categorías. Más adelante veremos con más detalle el uso que se puede hacer de tales valores.

Como destaca Zeisel en su clásico libro *Digalo con números* (1974, pág. 38), el contenido estadístico de las tres formas de calcular los porcentajes en una tabla es el mismo, pero al calcular los porcentajes en diferentes sentidos se acentúan distintas distribuciones y se ofrecen diferentes comparaciones. Así, del examen de la tabla 7.5.a se puede concluir que los varones se identifican preferentemente con posiciones ideológicas de izquierda y centro —el 42,9 por 100 se identifica con la izquierda y otro tanto lo hace con el centro, mientras que sólo un 14,2 por 100 lo hace con la derecha—, y las mujeres se identifican en mayor proporción con las posiciones de derecha y, en menor grado, con el centro e izquierda —en concreto, el 50 por 100 se identifica con la derecha, y sólo el 25 por 100 lo hace con la izquierda y centro, respectivamente.

Conclusiones de otro tipo se alcanzarán si examinamos la tabla 7.5.b. Así, si nos fijamos en la categoría ideológica del centro, se observa que hay más varones, el 60 por 100, que mujeres, el 40 por 100, entre los que se identifican con tales posiciones. De igual modo se pueden analizar las categorías de izquierda y derecha, comprobando para cada una de ellas la composición de varones y mujeres que con ellas se identifican.

Por último, cabe observar la distribución de los datos que se presentan en la tabla 7.5.c. Aquí, los números de cada celdilla representan los porcentajes de casos en relación al total N . La lectura de esta tabla nos permite, por ejemplo, concluir que la categoría más numerosa es la de mujeres de derechas, el 26,7 por 100, mientras que la categoría menos numerosa es la de varones de derechas, el 6,7 por 100; o también se puede observar que existe el mismo porcentaje de varones de izquierdas que de derechas, el 20 por 100, y lo mismo ocurre con las mujeres, entre las que un 13,3 por 100 se manifiesta de izquierdas y otro tanto como de centro.

No siempre es posible determinar qué variable es la independiente (o «causa») y qué variable es la dependiente (o «efecto»). En tal caso no es aplicable la regla causa y efecto o variable independiente a variable dependiente, ya que las tablas que se obtengan al calcular los porcentajes en un sentido u otro ofrecerán diferentes interpretaciones. Como ejemplo vamos a considerar los siguientes resultados, obtenidos en una encuesta sobre temas de actualidad realizada por el Centro de Investigaciones Sociológicas entre la población española en febrero de 1981. Una de las preguntas incluidas en el cuestionario hacía referencia a la per-

TABLA 7.6

percpción de la situación económica general del país segun las preferencias ideológicas de la población

Situation econ.

7.6.c. Porcentajes calculados en el sentido de las filas.

Total	100	60	36	4
Debrechha	100	62	36	4
Centro	100	33	43	4
Axum	100	2	4

Con frecuencia, los sociólogos preparan tablas que, sin dejar de con-
tener relaciones bivariables, son más complejas que las que hemos te-
nido ocasión de ver en las páginas anteriores. Así, se pueden confeccio-

1.2.b. Distribuciones condicionales más complejas

Otras consideraciones se podrían realizar al final de la comparación de las distribuciones condicionales que se incluyen en la tabla 7.6, sobre todo si situamos tales resultados en el contexto de la estructura social de la sociedad española. Pero basten las breves consideraciones aquí formuladas para poner de manifiesto las posibilidades de este tipo simple, pero eficaz, de análisis estadístico.

La lectura de la tabla 7.6c ofrece resultados diferentes a los anteriores. Así, por ejemplo, se puede observar que, tanto entre los que se identifican con la izquierda como entre los que se identifican con la derecha, las distinciones de los que se identifican la situación económica como buena, regular o mala son muy semisiguientes, mientras que entre los que se identifican con el centro político es mayor, relativamente, la proporción de los que consideran la situación regular que en el resto, mientras que disminuye la proporción de los que consideran la situación económica mala.

En la tabla 7.6 se presentan las frecuencias absolutas que se obtiene-
n en el sentido de la preferencia «ideológica». A partir de tales resultados se han
calculado los porcentajes en el sentido de las columnas (tabla 7.6.b) y
la variable «preferencia ideológica». Al cruzar la situación económica por
el sentido de la preferencia ideológica se obtienen los resultados
que se presentan en la tabla 7.6.

FUENTE: Resultados parciales de la Tabla 6 correspondiente al Barómetro de Opinión Pública, febrero 1981, RE/S, núm. 15, 1981, pág. 189.

Total	100	60	36	4
Derecha	100	62	36	2
Centro	100	33	43	4
Izquierda	100	33	43	4

nar tablas en las que las distribuciones condicionales hacen referencia a más de una variable independiente, y en las que la atención se concentra en la observación de la categoría o categorías más relevantes de la variable dependiente. La tabla 7.7 ilustra con claridad lo que venimos diciendo.

TABLA 7.7

*Perfil del español interesado y desinteresado por la política
(Porcentajes de encuestas nacionales de 1976 y 1980)*

Variables socioeconómicas	Con mayor interés		Con menor interés	
	1976	1980	1976	1980
<i>Sexo:</i>				
% de hombres ...	43	43	33	27
% de mujeres ...	23	23	54	49
<i>Edad:</i>				
% menores de 25 años ...	43	47	29	27
% más de 64 años ...	19	13	62	57
<i>Ocupación:</i>				
Más cualificados ...	51	61	28	16
Menos cualificados ...	16	14	65	55
<i>Municipio:</i>				
Más de 500.000 habitantes ...	41	43	35	38
Menos de 2.000 habitantes ...	30	19	44	52

FUENTE: LÓPEZ PINTOR, Rafael: «El estado de la opinión pública», REIS, 13, 1981, pág. 30.

En un estudio sobre el interés por la política entre la población española, Rafael López Pintor trata de delimitar el perfil de los españoles interesados por la política y el de aquellos que no lo están. Para ello toma datos de dos encuestas sobre actitudes políticas de los españoles realizadas, respectivamente, en 1976 y 1980, y con ellos prepara la tabla 7.7. En esta tabla sólo se contienen datos referentes a los porcentajes de los individuos que manifiestan tener un interés alto y un interés bajo por la política, según el sexo, edad, ocupación y municipio de residencia de los entrevistados. A su vez, de estas cuatro variables independientes sólo se incluyen en la tabla las dos categorías que aparecen más relevantes para el análisis buscado. Naturalmente, el sexo aparece dicotomizado en «hombres» y «mujeres», pero de la variable edad sólo se destacan los «menores de 25 años» y los «mayores de 64 años»; de la variable ocupación, los «más cualificados» y los «menos cualifi-

cados»; mientras que, en lo referente al municipio de residencia, sólo se incluyen los que residen en municipios de más de 500.000 habitantes y los que lo hacen en municipios menores de 2.000 habitantes.

De este modo, se destacan tan sólo los valores extremos o perfiles más acusados de los interesados o no por la política. Así, de la observación de la tabla anterior cabe concluir que el perfil tipo del español más interesado por la política es el de un hombre, de edad menor de veinticuatro años, con ocupación cualificada y residente en una gran ciudad. El perfil de la persona más desinteresada por la política sería el de una mujer, de edad superior a los sesenta y cuatro años, de baja cualificación ocupacional y residente en una zona rural. En una sola tabla, pues, se ha podido condensar un gran volumen de información, que ha permitido obtener conclusiones muy generales. Ni que decir tiene que este tipo de tabla es de gran interés analítico, sobre todo en aquellos estudios, como los realizados a través de encuesta, que permiten obtener un gran volumen de información que, necesariamente, se ha de resumir para poder alcanzar resultados globales. Obsérvese igualmente que los porcentajes no suman 100, ya que sólo se han incluido unas pocas categorías, y éstas, en consecuencia, no son exhaustivas.

7.3. CARACTERÍSTICAS DE UNA ASOCIACIÓN DE DOS VARIABLES

Cuando establecemos una clasificación cruzada de dos variables, nuestro interés se centra sobre todo en el conocimiento de la forma en que se distribuye la variable dependiente para las diferentes categorías de la variable independiente o causal. La forma en que se relacionan dos variables se denomina *asociación entre dos variables*. Volviendo a la tabla 7.3, se observa que a medida que disminuye el nivel de religiosidad se incrementa la proporción de personas que se identifican con las posiciones ideológicas de izquierda y, viceversa, a medida que aumenta el nivel de religiosidad se incrementa la proporción de personas identificadas con posiciones de derechas. Ese es, pues, el tipo de asociación que cabe observar entre las variables religiosidad e ideología.

Tal como señalan Loether y McTavish (*op. cit.*, 185), de igual modo que, al estudiar las distribuciones univariadas, éstas quedaban caracterizadas mediante el estudio de su tendencia central, variación o dispersión y forma, asimismo, se puede caracterizar la relación entre dos variables mediante el estudio de las siguientes características: 1) *existencia o no de una asociación*; 2) *la fuerza de la asociación*; 3) *la dirección de la asociación*, y 4) *la naturaleza de la asociación*.

A continuación vamos a estudiar con cierto detalle las cuatro características y, más adelante, se desarrollarán algunos índices que se pueden utilizar para medirlas.

Ya hemos dicho anteriormente que existe una asociación entre dos

Por su parte para ulteriores cálculos, vamos a ver a continuación cómo se calcula una tabla o modelo de no sociabilidad de una tabacalera variable cada quiebra. El problema consiste en calcular las frecuencias de cada celdilla a partir de los datos totales, de forma que la distribución de los datos no ofrezcan asociaión alguna. Supongamos que partimos de la siguiente distribución de frecuencias absolutas entre los variables dicotómicas:

estadística descriptiva bivariante 219

(X)	I	II	Total	(Y)
				Total
II	a	b	52	I
III	c	d	40	II
			92	Total

Para obtener los valores de a , b , c y d de forma que no exista rela-
ción entre las variables (X) e (Y) hay que aplicar el siguiente razo-

En general, la frecuencia esperada /, para una determinada celdilla se calcula multiplicando el total de la fila correspondiente a la celdilla por el número global N de casos. El dividendo es el producto obtenido que corresponde igualmente a dicha celdilla, divisor es total de la columna que corresponde por el número global N de casos.

$$\frac{N}{m \cdot m} = t_{12}$$

Aplicando la fórmula [7.1] a los datos que se continúen en la tabla anterior, se calculan de immediao los valores de a , b , c y d :

Columna 1	Columna 2	Columna 3	Columna 4
1	2	3	4
2	3	4	5
3	4	5	6
4	5	6	7

Y, es el total para la columna 1 de la tabla; n , es el total para la fila 1; N , es el número total de casos.

$$a = \frac{52 \cdot 30}{92} = 16,9$$

$$b = \frac{52 \cdot 62}{92} = 35,0$$

$$c = \frac{40 \cdot 30}{92} = 13,0$$

$$d = \frac{40 \cdot 62}{92} = 27,0$$

$$c = \frac{92}{40 - 30} = 13.0 \quad d = \frac{92}{40 - 62} = 27.0$$

Otra forma de decir si existe o no asociación entre dos variables consiste en comparar las frecuencias observadas en la tabla con las frecuencias que cabría esperar si no existiera asociación, o *frecuencias esperadas*. Si al comparar la tabla de datos reales con la tabla de no asociación no se observa diferencia alguna, cabe hablar entonces de que no existe una asociación.

Anteriormente se combinó de las comparaciones se pude en rea- lizar mediana el cálculo de los espesores, que son las diferencias por centrales calculadas en la dirección en que se han realizado los porcen- tajes. Así, para la categoría «poco» metros por la política, $= 29 - 19 = 10$ por 100. Pues bien, cuando hay asociación entre dos variables, la mayor parte de los espesores al calcular las diferentes categorías son diferentes de cero, mientras que cuando todos los espesores son cero no existe asociación alguna entre las variables. La idea de ausencia de aso-

FUENTE: Barómetro de Opinión Pública, Sep. 1981, REIS, núm. 16, 1981, pág. 224.

	Mujeres	Vaginas	Mujeres	Total
■ Muchos	7	4	11	100
■ Regular	23	19	27	100
■ Poco	24	19	32	100
■ Nadie	43	54	32	100
■ NS/NC	3	4	2	100
Total	100	100	100	(1.1193)

arribables cuando la distribución de una variable difiere de algún modo entre las divisiones categorías de la segunda variable. En su forma más general, y una vez calculados los porcentajes de la forma apropiada, se sucede decir que existe una asociación entre dos variables cuando las correspondientes distribuciones condicionales son diferentes entre sí. Así, por ejemplo, en la siguiente tabla, que muestra la variable «sexo» con la variable «interés por la política», se puede afirmar que existe una asociación entre ambas:

Una vez calculados los valores esperados en cada casilla, se pueden comparar los valores observados f_{ij} , o reales, de la tabla con los valores esperados f_{eij} . La comparación se realiza restando el valor esperado de cada celdilla del valor observado de la celdilla correspondiente. Este valor, denominado $\delta = f_{ij} - f_{eij}$, se calcula para cada celdilla de la tabla. Mientras algunos de los valores δ obtenidos sean diferentes de cero, se puede hablar de la existencia de algún tipo de asociación entre las dos variables. Si todos los deltas son cero, entonces se puede afirmar que no existe asociación alguna entre las variables, o, dicho con otras palabras, existe independencia estadística entre las dos variables.

Ahora bien, no es lo mismo, desde el punto de vista de la asociación entre dos variables, que los valores ϵ o δ sean altos o bajos. Aquí conviene introducir la noción del grado o fuerza de la asociación entre dos variables. Cuando los valores ϵ o δ son elevados cabe hablar de un alto grado de asociación o de una fuerte asociación entre las variables, mientras que si tales valores son pequeños se trata de una débil asociación o de un bajo grado de asociación.

Existe un problema con la utilización de los valores de ϵ y δ , y es que resulta difícil determinar con precisión el significado de un valor determinado, aparte de revelar la existencia o no de una asociación, ya que no existe una escala con un valor mínimo y un valor máximo entre los que puedan variar los valores obtenidos de ϵ y δ . Por dicha razón se utilizan con mayor frecuencia otro tipo de índices «estandarizados» que varían, de una forma fija, predeterminada e interpretable, entre un valor mínimo de no asociación y un valor máximo de mayor asociación. Más adelante estudiaremos los índices o coeficientes estandarizados de mayor uso en la investigación empírica en sociología.

Por lo que se refiere a la tercera de las características enunciadas, la dirección de la asociación, sólo cabe hablar de ella cuando las variables se han medido, como mínimo, al nivel ordinal. Con variables nominales o clasificadoras no cabe hablar de dirección de la asociación. Cuando, en una tabla, la tendencia de variación conjunta de las dos variables es a que los valores altos de una variable se correspondan con los valores altos de la segunda variable (y los valores bajos se corresponden igualmente), cabe hablar de la existencia de una asociación positiva. Así, en el siguiente ejemplo, con datos ficticios entre el nivel de ingresos y el nivel de satisfacción general, la dirección de la asociación entre ambas variables es positiva, ya que, a mayor nivel de ingresos, más elevado es el nivel de satisfacción:

Nivel de satisfacción	Nivel de ingresos		
	Bajo	Medio	Alto
— Bajo	60	40	30
— Medio	30	40	45
— Alto	10	20	25

Por el contrario, cuando los valores superiores de una variable se corresponden con los valores bajos de la segunda, y los valores altos de ésta se corresponden con los valores bajos de aquélla, se dice entonces que la dirección de la asociación es negativa. Así, por ejemplo, al estudiar la relación existente entre el nivel de ingresos de los individuos y el grado de anomia que padecen se observa una asociación negativa, ya que los individuos de ingresos altos tienden a tener un grado menor de anomia que los individuos de ingresos bajos, que padecen un mayor grado de anomia, como se puede observar en el siguiente cuadro:

Grado de anomia	Nivel de ingresos		
	Bajo	Medio	Alto
— Bajo	20	40	55
— Medio	30	25	45
— Alto	50	35	20

Finalmente, nos queda por analizar la cuarta característica de la asociación entre dos variables. La naturaleza de una asociación se refiere a la forma general en que se distribuyen los datos en la tabla. Habitualmente, dicha forma general o modelo se describe mediante el examen de las distribuciones de los porcentajes. En unos casos la distribución es irregular, distribuyéndose las diferencias elevadas o las aproximaciones entre cada par de porcentajes de una manera desigual, mientras que en otros casos se produce una progresión uniforme de las diferencias porcentuales desde las categorías bajas a las altas de las variables. Cuando, al pasar de una categoría a otra de una variable, el número de casos tiende a incrementarse (o disminuir) de una forma bastante homogénea entre las correspondientes categorías de la otra variable, se produce una asociación «lineal», esto es, que los casos se concentran en la variable dependiente siguiendo una línea recta. Las asociaciones lineales simples tienen un gran valor en la estadística en general, y en la investigación sociológica en particular, como modelos de asociaciones simples, aunque con frecuencia los datos sociológicos se distribuyen siguiendo formas

Por todo ello, el coeficiente de chi-cuadrado no se utiliza como medida de asociación, aunque, como se ha dicho anteriormente, si se amplia el rango de la estadística inferencial. Otros coeficientes basados en chi-cuadrado tratan de aprovecharse de las ventajas que ofrece dicho coeficiente, a la vez que tratan de superar, mediante el uso de estimaciones, la limitación que impone la hipótesis nula.

El coeficiente chi-cuadrado es siempre un número positivo, y se hace cero en las tablas en las que no hay asociación entre las variables. Si en embargo, el límite superior del coeficiente χ^2 no es fijo, sino que vale $N(K-1)$, en donde N es el tamaño de la muestra y K es el número de filas o columnas en la tabla, segün sean unas u otras las que representan el número más pequeño. Para una tabla 2×2 , el límite superior de la magnitud de χ^2 es N . Por tanto, dados dos tablas que tengan una asociación idéntica en su forma porcentual, si una de ellas se basa en un número doble de casos que en la otra, su valor de χ^2 será el doble que en la tabla basada en el menor número de casos.

Este índice se utiliza más en la estadística inferencial, para la prueba de hipótesis, que en la estadística descriptiva, para medir el grado de asociación entre dos variables, ya que presenta ciertos problemas al tratar de estandarizar sus valores. Sin embargo, al tener una «distinción libre», se convierte en una prueba muy útil para variables no-
nadas o categóricas.

$$e^2 = \sum_{\lambda} \left(\frac{f_\lambda}{A_\lambda} \right)^2$$

es igual a la suma de todas las diferencias que se pudieren establecer entre los valores observados y esperados. Ahora bien, este índice es muy deficitario, ya que depende, en primer lugar, del tamaño de los valores esperados y, además, los valores de los errores individuales se pierden al calcular el cuadrado por el problema de los signos) y se suman sus cuadrados (con lo que desaparecen los diferentes detalles, se pierden en buena medida si, en lugar de solo los efectos distorsionantes que se producen en cierto modo los efectos de los errores de casos. Los valores así obtenidos dan lugar a una medida nómica de los errores de los datos que no depende de la asociación de los datos, sino de las condiciones especiales que cumplir los datos.

En el caso anterior ya estudiamos la prueba del chi-cuadrado para una sola variable, que se basa en los cálculos que acabamos de describir, de igual manera que encotramos el chi-cuadrado cuando se realizan las pruebas de decisión estadística con dos o más variables. Pero ahora volvemos a la medida de asociación que acabamos de estudiar. La medida resultante de sumar todos los coeficientes anteriores para cada celdilla se denomina chi-cuadrado (χ^2):

Las medidas estandarizadas o tipificadas de asociación suelen ser simples proporciones o cocientes (ratios) que son sensibles a los cambios que se producen en el grado de asociación y, en algunos casos, en la dirección y naturaleza de la misma. De lo que se trata es de conseguir indicadores que reflejen realmente la variación de los aspectos de las variables cruzadas y que sean más sensibles a las variaciones de las variables cruzadas para la asociación, como pueden ser el número de filas y columnas o el número total de casos en que

Una medida de asociación del tipo que acabamos de describir es una medida de asociación estandarizada o tipificada, ya que los valores respectivos obtendrán la misma magnitud en diferentes tablas se pue- den comparar entre sí. Así, por ejemplo, si al cruzar la variable «inte- reses por la política» con la edad se obtiene una asociación de +0,52, y al hacerlo con la variable nivel de ingresos se obtiene una asociación de +0,35, podemos afirmar que la asociación entre una variable y el nivel de ingresos es más fuerte que la asociación entre una variable y el sexo.

Tal como se ha indicado anteriormente, el investigador necesita disponer de medidas que en un solo índice indiquen la existencia, grado y dirección de la asociación entre dos variables. Habitualmente, lo que se busca es una medida cuyos valores puedan variar a lo largo de una escala desde un valor mínimo, que indique una relación negativa, hasta un valor máximo, que indica una asociación positiva, pasando por el centro de la asociación.

INDEPENDENCIA ESTADÍSTICA Y ASOCIACIÓN PERFECTA

74. LA OBTENCIÓN DE MEDIDAS DE ASOCIACIÓN ENTRE DOS VARIABLES:

Cuávillenes o de otra naturaleza. Más adelante volveremos con mayor detalle a tratar este tema.

terminadas correcciones, sus deficiencias o limitaciones. Así, se puede utilizar el coeficiente de «contingencia cuadrática media» o *fi-cuadrado*, ϕ^2 , que se define simplemente como el valor de chi-cuadrado dividido por N :

$$\text{fi-cuadrado: } \phi^2 = \frac{x^2}{N} \quad \circ \quad \phi = \sqrt{\frac{x^2}{N}} \quad [7.3]$$

El valor de ϕ varía entre 0 —para el caso de independencia estadística— a un máximo de +1 —cuando existe una asociación perfecta—, en cualquier tabla de tamaño $2 \times K$, pudiéndose interpretar su magnitud como una medida del grado de asociación. Sin embargo, presenta el inconveniente de que, en tablas que contengan más de dos categorías en cada variable, el valor máximo de ϕ sobrepasa la unidad, dado que el límite superior de $x^2/N(K-1)$, se convierte en tal caso en un valor superior a N . El valor máximo de $\phi^2 = K-1$, en donde K representa el número más pequeño, bien de las filas o bien de las columnas.

El propio inventor del chi-cuadrado, el inglés Karl Pearson (1857-1936), considerado por muchos como el auténtico fundador de la estadística moderna, suministró una solución parcial a las anteriores limitaciones, mediante el desarrollo del «coeficiente de contingencia» o *coeficiente C de Pearson*. La fórmula para C es la siguiente:

Coeficiente de
contingencia

$$C = \sqrt{\frac{x^2}{x^2 + N}} \quad [7.4]$$

El coeficiente C no puede ser superior a la unidad, con independencia del tamaño de la tabla, ya que el coeficiente x^2 aparece tanto en el denominador como en el numerador, y aquél es siempre mayor que éste, ya que contiene la suma $x^2 + N$, que será siempre superior a x^2 . En su valor mínimo, el coeficiente C puede llegar a ser cero cuando, en los casos de ausencia de asociación, el valor de x^2 sea también cero, pero nunca alcanza exactamente la unidad, aunque hubiera asociación perfecta, por la razón anteriormente apuntada de que el denominador es siempre superior al numerador en la expresión [7.4]. Para una tabla cuadrada, es decir, una tabla en la que el número de filas sea igual al número de columnas, el valor máximo de C se puede calcular a partir de la expresión siguiente:

$$C \text{ máximo} = \sqrt{\frac{K-1}{K}}$$

Para tablas
cuadradas (= nº de filas
y columnas).

en donde K es el número de filas (o de columnas) en una tabla cuadrada. Así, por ejemplo, para una tabla 2×2 , el valor máximo de C es 0,707;

para una tabla 4×4 , el valor máximo de C es 0,87, y para una tabla 5×5 , el C máximo es 0,89. Así, pues, utilizando el coeficiente C no se pueden realizar comparaciones con esta medida de asociación entre tablas de diferentes tamaños.

Otros autores han tratado de mejorar la obtención de un coeficiente de asociación que pueda utilizarse para comparar tablas de diferente tamaño, es decir, que se pueda disponer de un coeficiente suficientemente estandarizado o normatizado. El *coeficiente T de Tschruprow* corrige el problema del límite superior de C mediante una ligera modificación del denominador de la expresión [7.4], de tal modo que incluya un valor que refleje el número de celdillas de la tabla. En otras palabras, se trata de introducir el concepto de los grados de libertad en la fórmula del coeficiente de asociación. Parece ser que el propio Pearson nunca llegó a comprender el concepto de grados de libertad en relación tanto con el chi-cuadrado como en relación con el cálculo de los errores de probabilidad (H. M. Walker, 1978, pág. 695). Por esa razón han tenido que ser otros autores los que se preocuparon de obtener medidas de asociación mejor normatizadas. Recordemos que, en una tabla de n filas y m columnas, los *grados de libertad** $df = (n-1) \cdot (m-1)$, es decir, es igual al número de filas menos uno multiplicado por el número de columnas menos uno.

Pues bien, el coeficiente T de Tschruprow se define del siguiente modo:

$$\text{Coeficiente } T. \quad T = \sqrt{\frac{x^2}{N(df)}} \quad [7.5]$$

El coeficiente T representa un avance en la búsqueda de una medida de asociación que esté adecuadamente estandarizada o normatizada para cualquier tipo de tabla. En efecto, el límite superior de T vale la unidad, con independencia del tamaño de la tabla, en tanto que ésta sea cuadrada, es decir, que el número de filas sea igual al número de columnas. Ahora bien, para tablas que no son cuadradas, el valor de T no puede alcanzar la unidad, aunque su valor máximo sea constante para tablas con idénticos grados de libertad.

Otro coeficiente, la *V de Cramer*, trata de resolver el problema de la estandarización o normatización mediante la sustitución en la ex-

* El concepto de grados de libertad lo hemos estudiado en el capítulo introductorio a la estadística inferencial. De una forma intuitiva, su concepto se puede entender en el estudio de las tablas bivariales, señalando que en una tabla 2×2 en donde $df = (2-1)(2-1) = 1$, se puede conocer una frecuencia esperada conociendo una frecuencia observada en una celdilla. El resto se puede calcular por substracción, ya que los marginales son fijos y, por tanto, conocidos. Se tiene, pues, un grado de libertad en elegir la frecuencia de una celdilla antes de que se puedan determinar las restantes frecuencias. En una tabla 3×3 , se han de elegir cuatro frecuencias de celdillas antes de determinar el resto, esto es, tiene cuatro grados de libertad, y así para otros tamaños de la tabla.

Estadística descriptiva bivariada 227

En una tabla 2×2 , estos signifICA que las células de una matriz diagonales contienen valores y las de la segunda diagonal tienen cero, como se observa a continuación:

que se corresponde con las mismas categorías de ambas variables, mientras que la tabla b) refleja una «asociación negativa», pues las que se observan en la tabla a) entre las categorías de ambas variables.

Vemos a través de un ejemplo hipotético el funcionamiento de es-
as concepciones. Supongamos que tratamos de contrastar la teoría de que
los delitos por consumo de drogas son más elevados en las grandes
ciudades en relación a las pequeñas ciudades y zonas rurales. Pues bien,
el modelo de la asociación preferida singulariza que, al distinguir los
tots en una tabla que cruzase la variable «freuencia de delitos por
consumo de drogas» por la variable «tamano del lugr de residencia»,
nos sumo de drogas» por la variable «tamano del lugr de residencia».
entras los detalles de tal tipo se concentran en las grandes ciudades,
es detiles que en las ciudades pequeñas no se produciría ninguna
desviación de la media. Cuadúller desviación a esta forma
de datos correspondientes a las dos variables significa una aso-
ción no perfecta.

urales que se observan en la figura 1, como se observa en la

		Allocation per unit			Allocation per unit		
		I	A	B	I	A	B
		II	0	0	II	0	0
		(X)	(Y)	(Z)	(X)	(Y)	(Z)
(X)							

Se dice que una tabla bivariante refleja una asociación perfecta cuando todos los casos de la tabla se concientran en una diagonal, lo que significa que cada valor de una variable se encuentra asociada con un solo valor de la segunda variable, de tal modo que para cualquier categoría de la variable independiente, solo será diferente de cero una categoría de la variable dependiente, mientras que el resto de las células

$$[7.6] \quad \frac{t \cdot N}{\epsilon^x} = A$$

presentación de T de los grados de libertad d_f por un valor ; que representa el número más pequeño de las dos cantidades, $n - 1$ o $m - 1$, siendo n el número de filas y columnas, respetivamente. Así, pues, la fórmula de la V de Cramér es como sigue:

Pues bien, el coeficiente Q de Yule se calcula a partir de los productos cruzados de las celdillas de una de las diagonales ad y de las celdillas de la segunda diagonal bc . El coeficiente Q de Yule se calcula mediante una fórmula como sigue:

$$\text{puede ser } Q \text{ de Yule.} \quad Q = \frac{ad - bc}{ad + bc} \quad [7.7]$$

c/u. normal.

Cuando la frecuencia de una de las celdillas sea cero, entonces el valor de Q es $+1,0$ o $-1,0$, según la dirección de la asociación. El coeficiente Q se puede utilizar con variables nominales y, cuando alcanza el valor de la unidad en una tabla 2×2 , refleja la existencia de una asociación perfecta.

7.4.1. Medidas simétricas y asimétricas de asociación

Finalmente, vamos a señalar una última distribución de las medidas de asociación que tiene interés para la investigación sociológica. Hay medidas de asociación que distinguen entre la variable independiente (o «causa») y la variable dependiente («efecto»), mientras que otras medidas de asociación no realizan tal distinción.

Pues bien, a las medidas de asociación que no distinguen entre variables independientes o dependientes se les denomina *medidas simétricas*. Tales medidas reflejan tan sólo la fuerza (y dirección) de la relación entre dos variables, y no distinguen entre los papeles asignados a cada variable. Los coeficientes vistos con anterioridad, tales como la Q de Yule, el coeficiente ϕ , la C de Pearson, la V de Cramer o la T de Tschruprow son ejemplos de medidas simétricas de asociación.

Por otro lado, hay medidas de asociación que requieren para su cálculo que se distinga previamente entre la variable independiente y la variable dependiente. Se trata de *medidas asimétricas de asociación*, que están orientadas, en general, a medir la capacidad e influencia de una variable independiente en la predicción de los valores de la variable dependiente. Buena parte de los coeficientes que vamos a estudiar en los próximos capítulos son de tipo asimétrico, aunque ya en este mismo capítulo hemos tenido ocasión de estudiar una de tales medidas. En efecto, el coeficiente ϵ , que, como se recordará, es una simple diferencia entre porcentajes, ofrece diferentes valores según sea el sentido en que se calculen los porcentajes, es decir, según sea una u otra la variable que se considera independiente. Naturalmente, al variar las bases sobre las que se calculan los porcentajes, así variarán los valores de ϵ . De todas maneras, este coeficiente apenas se utiliza en la práctica de la investigación sociológica, porque no se trata de una medida normatizada, como las que veremos en el próximo capítulo.

7.5. TERMINOLOGÍA

Se recomienda la memorización y comprensión del significado de cada uno de los términos y conceptos siguientes:

- Distribución porcentual bivariable.
- Distribución condicional.
- Asociación entre dos variables:
 - Existencia de la asociación.
 - Fuerza o grado de la asociación.
 - Dirección de la asociación.
 - Naturaleza de la asociación.
- Frecuencias observadas.
- Frecuencias esperadas.
- Independencia estadística.
- Asociación perfecta.
- Asociación positiva.
- Asociación negativa.
- Coeficiente ϵ .
- Coeficiente δ .
- Coeficiente ϕ .
- Coeficiente chi-cuadrado.
- Coeficiente C de Pearson.
- Coeficiente T de Tschruprow.
- Coeficiente V de Cramer.
- Coeficiente Q de Yule.
- Grados de libertad.
- Medidas simétricas de asociación.
- Medidas asimétricas de asociación.

EJERCICIOS

1. De los siguientes pares de variables: ¿Cuáles están formados por variables independientes entre sí (es decir, no es posible *a priori* especificar una ordenación causal o temporal entre ellas)?; ¿cuáles están formados por variables que están relacionadas entre sí condicionalmente? Para estos últimos pares de variables, especificar para cada par qué variable, desde un punto de vista lógico, antecedente de la otra.
 - 1) Tamaño de familia y religiosidad de los cónyuges.
 - 2) Edad y región de nacimiento.
 - 3) Interés por la política y nivel de educación.

5. En una encuesta sobre actitudes de la población hacia el aborto, las opiniones sobre la legalización o prohibición del aborto se distribuyeron del siguiente modo, teniendo en cuenta la identificación ideológica de la población:

Ideología	¿Debe permitirse el aborto?		
	En ningún caso	Por necesidad	Por decisión libre
Izquierda	100	280	370
Centro	250	410	90
Derecha	370	280	60

A la vista de la anterior distribución, ¿se puede afirmar que existe asociación entre ambas variables? ¿De qué tipo es? En caso afirmativo, calcular el grado de asociación mediante el coeficiente de contingencia C de Pearson. Comparar el valor obtenido con el valor máximo de C que se podría obtener para una tabla del tamaño como la presente.

BIBLIOGRAFIA

- JIMÉNEZ BLANCO, J., et al.: *La Conciencia Regional en España*, Madrid, CIS, 1977.
 LOETHER, H. J. y D. G. McTAVISH: *Descriptive Statistics for Sociologists*, Boston, Allyn & Bacon, 1974.
 LÓPEZ PINTOR, R.: «El estado de la opinión pública y la transición a la democracia», *REIS*, núm. 13, 1981, págs. 7-47.
 WALKER, Helen W.: «Karl Pearson», en W. H. Kruskal y J. M. Tanur (eds.) *International Encyclopedia of Statistics*, New York, Free Press, 1978, págs. 691-698.
 ZEISEL, Hans: *Digalo con números*, México, Fondo de Cultura Económica, 1974. (e. o. 1947).

Capítulo 8

MEDIDAS DE ASOCIACIÓN PARA VARIABLES NOMINALES Y ORDINALES

Son muy variadas las medidas de asociación de que puede disponer un sociólogo interesado en el estudio de relaciones bivariadas. En el capítulo anterior tuvimos ocasión de estudiar algunas de ellas basadas en el valor de delta, o diferencia entre la frecuencia observada y la frecuencia esperada. Pero algunos de los coeficientes estudiados en dicho capítulo no son de interés para el investigador social, ya que no están «normatizados» y, por lo tanto, no está recomendada su utilización comparativa entre diferentes tablas, y menos aún la interpretación del carácter de la asociación. En el presente capítulo vamos a estudiar las medidas de asociación basadas en el criterio de «reducción proporcional del error», por ser las más utilizadas por los sociólogos, y ello para las relaciones entre variables medidas a nivel nominal y a nivel ordinal. En el próximo capítulo continuaremos con el estudio de las medidas basadas en el mismo criterio de reducción del error, pero para el caso de variables de intervalo, con lo que abordaremos uno de los temas centrales de la estadística, el estudio de la regresión simple.

Dado el carácter introductorio del presente libro, no vamos a estudiar las medidas de asociación apropiadas para situaciones especiales, porque esperamos que, con el bagaje de técnicas estadísticas que se presentan aquí, el estudiante de sociología puede pasar a realizar por sí mismo una investigación empírica sólida. Por ello remitimos al lector interesado en medidas de asociación especiales a otros libros, tales como el de Freeman (1971), y algunos otros trabajos que se citan en la bibliografía, para que pueda estudiar y conocer las mismas.

8.1. MEDIDAS DE ASOCIACIÓN BASADAS EN EL CRITERIO DE «REDUCCIÓN PROPORCIONAL DEL ERROR» (RPE)

Un simple repaso al estudio de las diferentes medidas de asociación disponibles para el estudio de datos pone rápidamente de manifiesto la dificultad de encontrar un principio lógico consistente que sea

