

HERMILDA SUSANA RONDÓN TRONCOSO*

FECHA DE RECEPCIÓN: 17 DE ENERO DE 2015
FECHA DE EVALUACIÓN: 23 DE FEBRERO DE 2015
FECHA DE ACEPTACIÓN: 26 DE MAYO DE 2015

SOBRE EL CRUCE ENTRE VARIABLES CATEGÓRICAS

Crosstabs between nominal variables

Sobre o cruzamento entre variáveis categóricas

* Licenciada en Matemáticas y Física, de la Universidad del Tolima, Colombia; especialista en Estadística, de la Universidad Nacional de Colombia, Bogotá; magíster en Tecnología Educativa, del Tecnológico de Monterrey, México. Profesor asistente, de la Escuela Colombiana de Ingeniería Julio Garavito, Bogotá. Coordinadora de probabilidad y de estadística, de la Escuela Colombiana de Ingeniería Julio Garavito. Correo electrónico: susana.rondon@escuelaing.edu.co



Cómo citar este artículo: Rondón Troncoso, H. S. (2015). Sobre el cruce entre variables categóricas. *Revista de Educación y Desarrollo Social*, 9(2), 74-85.

RESUMEN

La mayoría de los libros tradicionales de estadística para pregrado abordan el tema sobre el cruce entre variables categóricas, extrayendo directamente la información que traen las tablas de valores observados y esperados para conseguir el valor puntual del estadístico de la Ji cuadrado. La información presentada de esta manera no le permite al estudiante comprender de dónde salen estos valores, ni tampoco apreciar el trabajo que hay al construir estas tablas, en especial para muestras grandes. El objetivo del presente artículo es dar a conocer una forma alternativa de realizar el cruce entre variables categóricas, y así permitir al estudiante realizar todo el proceso de construcción de las tablas de

valores observados y esperados. La metodología propuesta permitirá diseñar el archivo de datos con la participación de los estudiantes para un grupo de variables categóricas, de las que se seleccionan sexo y estado civil. Los valores de estas dos variables presentados por columnas se van colocando en una tabla en la casilla que corresponda hasta completar la muestra; esta tabla generará las tablas de valores observados y esperados que permitirán el cálculo del estadístico de prueba. El resultado se compara con el obtenido al usar una macro construida especialmente para este propósito.

Palabras claves: independencia entre variables, tablas de doble entrada, prueba Ji cuadrado, variables categóricas.

ABSTRACT

When it comes to working with crosstabs between nominal variables most of undergraduate statistics textbooks approach this topic by extracting the information found in observed and expected counts to get the exact chi-square value. Consequently, students can neither comprehend where these values come from, nor appreciate the worth of creating these tables, especially in larger samples. The primary objective of this paper is to propose an alternative to make crosstabs with nominal variables, allowing the student to develop the entire table building process for both observed and expected counts. The proposed methodology will allow to design the data file with the participation of the students for a group of nominal variables, of which two were chosen, Sex and Marital Status. The values of these two variables, shown in columns, are put in a table in their corresponding cell until the sample is complete. This table will generate the observed and expected counts used to calculate the test value. This result is compared with the one obtained using a macro programmed for this specific purpose.

Keywords: Variable independence, crosstabs, Chi-square test, nominal variables.

RESUMO

A maioria dos livros tradicionais de estatística para pre-graduação abordam a questão sob o cruzamento entre variáveis categóricas, extraindo diretamente a informação que trazem as tabelas de valores observados e esperados para conseguir o valor pontual do estatístico da Chi-quadrado. A informação apresentada desta forma não lhe permite ao aluno compreender de onde vêm esses valores,

nem apreciar o trabalho que precisa para construir essas tabelas, especialmente para grandes amostras. O objetivo deste trabalho é apresentar uma forma alternativa de fazer o cruzamento entre variáveis categóricas, e assim permitir que o aluno possa realizar todo o processo de construção de tabelas de valores observados e esperados. A metodologia proposta permitirá projetar o arquivo de dados com a participação dos estudantes para um grupo de variáveis categóricas, das quais sexo e estado civil são selecionados. Os valores destas duas variáveis apresentadas em colunas vão-se colocando em uma tabela no escaninho apropriado até completar a amostra; esta tabela vai gerar as tabelas de valores observados e esperados que vá a permitir o cálculo da estatística de teste. O resultado é comparado com o obtido quando se usa uma macro construída especialmente para este fim.

Palavras-chave: independência entre variáveis, tabelas bidirecionais, teste Chi-quadrado, variáveis categóricas.

INTRODUCCIÓN

Tradicionalmente, en estadística el tema sobre la prueba Ji cuadrado de independencia, correspondiente al análisis de encuestas, en libros de autores como Mendenhall (2010); Montgomery (2004); Sheldon, M. R. (2000); Devore, J. (2012); Navidi (2006); Seymor, L. (1991); Walpole (2012); Levin, R. R. (1996); Anderson (2008); Freund, J. E. y Garay, A. S. (1992); Lind, D. M., William, M. y Robert, D. M. (2004); Nieves, A. D. y Federico, C. D. (2010); Kazmier, L. y Díaz, M. A. (1993), entre otros, se presentan las tablas de valores observados y en

algunos casos también la de valores esperados ya elaboradas, de donde se toma la información que luego es reemplazada en el algoritmo de la prueba Ji cuadrado de independencia, para lograr así el cálculo del mismo. La metodología presentada así tiene el inconveniente de no permitirle al estudiante detectar el origen de los valores que vienen en las tablas mencionadas. Para el estudiante que recién comienza a aprender sobre este tema no es fácil apropiarse del concepto con solo tomar valores de unas tablas que luego reemplazará en una fórmula.

En un análisis de encuesta, lo que comúnmente se realiza son tablas de frecuencias, diagramas de torta y de barras, y se deja a un lado el cruce entre variables categóricas que le da mayor peso al análisis. Otro aspecto importante en la encuesta es el análisis de correspondencia para la clasificación de grupos de individuos según lo comentado por Clavijo (2005), pero esto requiere un nivel más alto de conocimientos estadísticos.

El tema central de este artículo es el cruce entre variables categóricas mediante la prueba Ji cuadrado de independencia. Años de experiencia docente han permitido observar que los estudiantes se apropian con mayor facilidad de este concepto para esta distribución cuando ellos mismos construyen las tablas de valores observados y esperados. Realizar el conteo de los datos de las dos variables que se seleccionaron para el cruce y que serán leídos de la tabla de datos o archivo de datos (tabla 2) logrará una mayor motivación, seguridad y apropiación del concepto sobre el tema en cuestión; asunto que se detallará en el apartado de procedimiento.

No obstante, se debe aclarar al estudiante que el cruce entre variables categóricas hace parte de la estadística descriptiva, que hasta

1900 y a pesar de sus limitaciones ha permitido hacer grandes aportaciones a la ciencia (Batañero y Godino, 2001).

El presente artículo proveniente de una experiencia pedagógica, propone calcular el estadístico de prueba de la Ji cuadrado de independencia paso a paso, lo que permite apreciar el proceso de construcción de las tablas. Para el ejercicio propuesto se tuvieron en cuenta treinta valores que fueron creados en la clase de estadística con la participación de los estudiantes y se pueden apreciar en la tabla 2. Posteriormente, se calculó el mismo estadístico con ayuda de una macro, que para este artículo se llamó (χ^2), y se demostró que los resultados son iguales. En caso de muestras grandes se podrá hacer uso de algún *software* estadístico como SSPS, Excel, SAS, entre otros.

PROCEDIMIENTO

El procedimiento que se detallará a continuación surgió de la experiencia con los estudiantes en el aula de clase, en la asignatura de Estadística, para el tema de análisis de encuestas, específicamente para el cruce entre variables categóricas. Generalmente, al estudiante se le proporciona un archivo ya elaborado, para que proceda a realizar su correspondiente análisis, pero cuando se tiene que enfrentar a ser él quien aplica la encuesta, no le es sencillo saber qué hacer con los datos logrados. De aquí nació la idea de que fueran ellos mismos los que en pequeños grupos de trabajo tuvieran que elaborar el archivo una vez aplicado el formulario, para poder finalmente proceder a realizar cruces entre variables categóricas y así poder ejecutar el análisis de una encuesta de forma más eficiente.

A continuación se darán algunas definiciones sobre tablas de contingencia. Una tabla de contingencia, según Arrondo (2014), se define como una “organización de filas y columnas, en cuyas casillas se expresa la frecuencia de ocasiones en las que se presenta el par valor_fila x valor_columna” (p. 2). Los autores Otero y Moral (2005) definen la tabla de contingencia como “una tabla de doble entrada, donde en cada casilla figurará el número de casos o individuos que poseen un nivel de uno de los factores o características analizadas y otro nivel del otro factor analizado” (p. 2).

Cuando se tiene una tabla de contingencia interesa ver si las variables representadas en las filas y columnas están relacionadas entre sí. En este caso, se está haciendo referencia a la asociación entre las dos variables, según lo comentado por Batanero y Díaz (2008). En general, una tabla de contingencia nos proporciona una forma resumida de representar datos de dos variables que se quieren estudiar, según Cañadas, Contreras, Arteaga y Gea (2013).

Sobre tablas de contingencia se ha venido trabajando desde muchos años atrás, según Inhelder y Piaget (1955) estos veían las dificultades

en la interpretación de las tablas de contingencia, y pensaban que “la comprensión de la asociación sería el último paso en el desarrollo del razonamiento sobre probabilidad” (p. 35).

A continuación se mostrará el procedimiento para realizar las tablas de contingencia, tanto de valores observados, como la de los valores relativos y la de valores esperados, con los datos del archivo que fueron logrados en el aula de clase. Aquí se podrá estudiar la relación entre el estado civil y sexo para un grupo de 30 personas. Las tablas se construyeron basándose en las ideas que tuvieron Inhelder y Piaget (1955), ya que como se comentó en el párrafo anterior desde esa época se elaboraron tablas similares a las propuestas para este artículo.

Empezaremos a construir la tabla de valores observados, para la cual se tendrá en cuenta un archivo que fue creado para este ejercicio (tabla 2). Del archivo se escogieron las variables sexo y estado civil (E. Civil); la tabla de doble entrada que se construirá llevará los nombres de estas dos variables y no importa la posición en que se coloquen en la tabla (donde aparece E. Civil, pudo haber estado sexo, o al contrario, tabla 1).

Tabla 1. Conteo de los 5 primeros valores observados

		ESTADO CIVIL			
		1=Soltero	2=Casado	3=Unión libre	4=Otro
SEXO	1=Masculino	11] ● ●	12] ●	13]	14]
	2=Femenino	21]	22] ●	23]	24] ●

Tabla 2. Archivo de datos inventado para la clase con (n=30)

IND.	EDAD	SEXO	E. CIVIL	GUSCAR	GUSMU	En el archivo hay 5 variables:
1	1	1	1	2	1	Edad categorizada: 1. Entre 17 y 20 años 2. Entre 21 y 25 años 3. Entre 26 y 30 años 4. Entre 31 y 40 años
2	2	2	2	1	2	
3	4	1	1	4	3	
4	3	2	4	2	2	
5	2	1	2	3	1	
6	1	1	3	4	2	
7	2	2	2	2	3	Sexo: 1. Masculino 2. Femenino
8	1	2	1	3	3	
9	2	1	2	1	2	
10	1	2	3	4	1	
11	3	1	2	4	1	Ecivil: Estado civil: 1. Soltero 2. Casado 3. Unión libre 4. Otro
12	4	2	1	1	1	
13	4	1	4	2	2	
14	3	1	2	3	3	
15	2	2	3	2	2	
16	1	2	2	4	1	
17	2	1	1	1	2	Guscar: Gustos por una carne 1. Lomo de res 2. Mojarra frita 3. Lomo de cerdo 4. Pollo con champiñón
18	3	1	2	2	3	
19	4	2	3	3	2	
20	2	1	2	2	1	
21	1	2	1	4	2	
22	2	1	2	2	3	
23	3	2	1	3	2	Gusmu: Gustos por 1 determinada música: 1. Pop 2. Rock 3. Electrónica
24	4	1	2	2	1	
25	2	2	3	4	2	
26	3	1	2	1	3	
27	1	2	3	2	2	
28	4	1	2	3	2	
29	2	2	1	1	3	
30	3	1	2	4	3	

Tabla 4. Valores observados y totales

		ESTADO CIVIL								TOTALES Marginales en Y
		1=Soltero		2=Casado		3=Unión libre		4=Otro		
SEXO	1=Masculino	11]	3	12]	11	13]	1	14]	1	16
	2=Femenino	21]	5	22]	3	23]	5	24]	1	14
TOTALES Marginales en X		8		14		6		2		30

Tabla 5. Valores observados logrados con la macro (Xi_Cuadrado)

	SOLTERO	CASADO	UNIÓN LIBRE	OTRO	TOTAL
Masculino	3	11	1	1	16
Femenino	5	3	5	1	14
TOTAL	8	14	6	2	30

A continuación se podrá apreciar la tabla de valores observados que se logró con la ayuda de la macro (Xi_Cuadrado). Como se puede apreciar, los resultados son los mismos que los logrados al realizar los cálculos a mano (tabla 5).

La siguiente expresión corresponde al algoritmo de la prueba de independencia de la Ji cuadrado (ecuación 1).

$$1. \quad \chi^{2*} = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i}$$

Cuando se tiene una tabla de contingencia se está interesado en ver si las variables representadas en las filas y columnas están relacionadas entre sí.

Ahora se podrá apreciar cómo quedan registrados los tres valores que se mencionaron en el anterior apartado (los tres hombres solteros, los once hombres casados y una mujer que se encuentra en otro estado de las categorías del estado civil). Posteriormente, se podrá observar la expresión para el estadístico de la prueba Ji cuadrado completa; esto por supuesto cuando se complete la tabla de valores esperados (ecuación 2).

$$2. \quad X^2 = \sum_i^n \frac{(O_i - E_j)^2}{E_j} = \frac{(3 - \quad)^2}{\quad} + \frac{(11 - \quad)^2}{\quad} + \dots + \frac{(1 - \quad)^2}{\quad}$$

La tabla de valores esperados se deriva de la tabla de valores relativos. A continuación se mostrará la tabla de valores relativos (tabla 6).

La construcción de la tabla de valores relativos es muy sencilla: teniendo en cuenta la tabla 4 se toman los 4 valores totales por fila (marginales en X: 8, 14, 6 y 2) y se multiplican por el primer valor total (16) de los marginales

en Y; cada producto se divide a su vez por el total de la muestra, que para este caso fueron 30. De igual forma, se procede a formar la segunda fila, pero en esta oportunidad los 3 valores (8, 14, 6 y 2) son multiplicados ahora por el segundo valor del total de los marginales en Y, es decir, 14. De igual forma, los cuatro productos de nuevo deberán ser divididos entre el tamaño de la muestra (30).

Al efectuar las operaciones de la tabla 6, se llega finalmente a la tabla de valores esperados (tabla 7).

A continuación se podrá apreciar la tabla de valores relativos lograda con ayuda de la macro (Xi_Cuadrado); se puede apreciar que los valores son los mismos que se lograron a mano y que están registrados en la tabla 7, pero con más cifras significativas (tabla 8).

Tabla 6. De valores relativos

$\frac{16 \times 8}{30}$	$\frac{16 \times 14}{30}$	$\frac{16 \times 6}{30}$	$\frac{16 \times 2}{30}$
$\frac{14 \times 8}{30}$	$\frac{14 \times 14}{30}$	$\frac{14 \times 6}{30}$	$\frac{14 \times 2}{30}$

Tabla 7. Valores esperados

4.27	7.47	3.2	1.07
3.73	6.53	2.8	0.93

Tabla 8. Valores esperados logrados con la macro (Xi_Cuadrado)

	SOLTERO	CASADO	UNIÓN LIBRE	OTRO
Masculino	4,26666667	7,46666667	3,2	1,06666667
Femenino	3,73333333	6,53333333	2,8	0,93333333

Ahora se podrán reemplazar estos valores en el estadístico de prueba como se muestra en la ecuación 3:

$$3. \quad X^2 = \sum_i^n \frac{(O_i - E_j)^2}{E_j} = \frac{(3 - 4.27)^2}{4.27} + \frac{(11 - 7.47)^2}{7.47} + \dots + \frac{(1 - 0.93)^2}{0.93} = 7.64$$

De esta forma, se podrá completar el algoritmo de la prueba de independencia Ji cuadrado. Solo se colocaron los dos primeros valores y el último como se hizo en la ecuación 2; lo importante es que a cada valor observado de la tabla 5 se le resta el correspondiente valor esperado de la tabla 7.

A continuación se puede apreciar el estadístico de prueba de la Ji cuadrado logrado con ayuda de la macro (Xi_Cuadrado); se observa que el resultado logrado fue de 7,63871173 (tabla 9) y es el mismo que el calculado en la ecuación 3 y que por aproximación arrojó un valor de 7,64.

Tabla 9. Ji cuadrado lograda con la macro

JI CUADRADO CALCULACO
7,63871173

Una vez logrados los cálculos con la ayuda de la macro (Xi_Cuadrado), es conveniente que el estudiante tenga en cuenta los siguientes pasos para completar todo el ejercicio.

Primer paso: se deben plantear las dos hipótesis, tanto la nula como la alterna, como sigue:

H₀: "Las variables son independientes"

H₁: "Las variables no son independientes"

Segundo paso: se calcula el estadístico de prueba ; este valor va en mayúscula y como ya se calculó se sabe con certeza que es = 7.64, valor logrado tanto a mano como con la ayuda de la macro (Xi_Cuadrado).

Tercer paso: se busca el correspondiente valor en tabla de la Ji cuadrado

$$X^2 (\gamma, \alpha) = X^2 (3, 0.05) = 7.815$$

El anterior valor se buscó teniendo en cuenta 3 grados de libertad. Los grados de libertad se calculan de acuerdo con el número de filas y el número de columnas de la tabla de valores observados. Como en las filas están los dos géneros (femenino y masculino), se realiza la diferencia entre las categorías del género que son 2 y 1 que es constante (2 - 1); en las columnas aparecen las cuatro categorías del estado civil, también se realiza esta diferencia, las cuatro categorías menos 1 que es constante (4 - 1), y quedan finalmente tres grados de libertad así:

Grados de libertad para la Ji cuadrado:

$$\gamma = (\text{Número de filas} - 1) (\text{Número de columnas} - 1)$$

$$\gamma = (2 - 1) (4 - 1) = 3 \text{ g.l. (grados de libertad)}$$

γ = Es una letra del alfabeto griego y se pronuncia Nú.

Ahora se busca el valor en la tabla de la Ji cuadrado, se tomará = 5 %; este valor se escoge a criterio del investigador.

$$X^2 (\gamma, \alpha) = X^2 (3, 0.05) = 7.815$$

Cuarto paso: se comparan los dos valores; el del estadístico de prueba y el valor encontrado en la tabla, de la siguiente manera:

Si $X^2 > x^2_{(\gamma, \alpha)}$, entonces se rechaza la hipótesis nula.

En este caso, como el valor que dio el estadístico de prueba fue de = 7.64, que resultó menor que el valor buscado en tabla que fue de: = 7.815, entonces no se rechaza la hipótesis nula. Esto significa que las variables son independientes, lo que indica que el estado civil no tiene nada que ver con ser hombre o ser mujer.

Es importante aclarar que la prueba Ji cuadrado de independencia tiene el inconveniente de no ser confiable para muestras pequeñas como la del ejemplo presentado para los 30 datos. Sin embargo, se hizo así por cuestión pedagógica y para facilidad del trabajo en clase y poder demostrar que por cualquiera de los dos métodos el resultado es el mismo (a mano y con ayuda de la macro Xi_Cuadrado).

CONCLUSIONES

El estadístico de la Ji cuadrado permite realizar pruebas de independencia, para variables categóricas, y brinda una mayor profundidad al análisis de encuestas. El cálculo para este estadístico se puede realizar a mano, siempre y cuando la

muestra sea pequeña por lo engorroso que resultaría el trabajo con una muestra grande. Los datos para este cálculo fueron tomados de un archivo que se elaboró con la información lograda a partir de una encuesta elaborada en clase. Este proceso logró que los estudiantes adquirieran destrezas y habilidades, y a su vez les facilitó la solución de ejercicios para tablas ya elaboradas y que vienen propuestas en los libros.

El cálculo a mano para el estadístico de la prueba Ji cuadrado de independencia le brinda al estudiante una mejor comprensión sobre el trabajo que tendría que realizar si hubiese elaborado una tabla para una muestra grande. De igual forma, sabrá con certeza que estos cálculos se pueden realizar con ayuda de algún *software* o de una macro (como la Xi_Cuadrado), lo cual facilitaría el trabajo y permitiría el análisis de forma más ágil y eficiente.

El estudiante debe tener claro que la prueba Ji cuadrado de independencia tiene el inconveniente de no ser muy confiable cuando la muestra es pequeña. La implementación de herramientas tecnológicas, en este caso un *software* o macro para un tema como el cruce entre variables categóricas, motiva al estudiante hacia el estudio de esta temática que tiene diversos usos y aplicaciones en la vida práctica.

REFERENCIAS

- ▶▶ Anderson, S. W. (2008). *Estadística para administración y economía* (10ma. edición). Cincinnati: Cengage Learning.
- ▶▶ Arrondo, V. M. (2014). Relaciones entre dos variables: una visión de urgencia. *Universidad de Sevilla*. Recuperado el 3 de marzo de 2015, de <http://asignatura.us.es/dadpsico/apuntes/RelacionesUrgencia.pdf>
- ▶▶ Batanero, C. y Godino, J. D. (2001). Análisis de datos y su didáctica. *Universidad de Granada*. Recuperado el 12 de enero de 2015, de <http://www.ugr.es/~batanero/pages/ARTICULOS/Apuntes.pdf>
- ▶▶ Bataeno, C. y Díaz, C. (2008). *Análisis de datos con Statgraphics*. Granada: Departamento de Didáctica de la Matemática.
- ▶▶ Clavijo, M.J. (2005). *Una introducción a la estadística general*. Ibagué: Universidad del Tolima.
- ▶▶ Cañadas, G. R., Contreras, J. M., Arteaga, P. y Gea, M. (2013). Problemática y recursos en la interpretación de las tablas de contingencia. *Revista Iberoamericana de Educación Matemática*, (34), 85-96. Recuperado el 14 de marzo de 2015, de <http://www.fisem.org/www/union/revistas/2013/34/archivo9.pdf>
- ▶▶ Devore, J. L. (2012). *Probabilidad y estadística para ingeniería y ciencias* (8va. edición). California: Cengage Learning.
- ▶▶ Freund, J. E. y Simon, G. A. (1992). *Estadística elemental* (8va. edición). México, D. F.: Prentice Hall.
- ▶▶ Inhelder, B. y Piaget, J. (1955). *De la logique de l'enfant à la logique de l'adolescent*. París: Presses Universitaires de France.
- ▶▶ Kazmier, L. (1993). *Estadística aplicada a la administración y a la economía* (2da.edición). Arizona: McGraw-Hill.
- ▶▶ Levin, R. R. (1996). *Estadística para administradores Schaum* (6ta. edición). México, D. F.: Prentice Hall.
- ▶▶ Lind, D. M., William, M. y Robert, D. M. (2004). *Estadística para administración y economía* (11va. edición). México, D. F.: Alfaomega.
- ▶▶ Mendenhall, W., Robert, J. B. y Barbara, M. B. (2010). *Introducción a la probabilidad estadística* (13va. Edición). México, D. F.: Cengage Learning.
- ▶▶ Montgomery, R. (2004). *Probabilidad y estadística aplicada a la ingeniería* (2da. edición). México, D. F.: Limusa Wiley.
- ▶▶ Navidi, W. (2006). *Estadística para ingenieros y científicos*. México, D. F.: McGraw-Hill.
- ▶▶ Nieves, A. D. y Federico, C. D. (2010). *Probabilidad y estadística para ingeniería un enfoque moderno* (1era. edición). México, D. F.: McGraw-Hill.
- ▶▶ Otero, J.V. y Moral, E. M. (2005). *Análisis de datos cualitativos*. Recuperado el 14 de marzo de 2015, de https://www.uam.es/personal_pdi/economicas/eva/pdf/tab_conting.pdf
- ▶▶ Seymour, L. (1991). *Probabilidad. Schaum* (1era. edición). México, D. F.: McGraw-Hill.
- ▶▶ Sheldon, M. R. (2000). *Probabilidad y estadística para ingenieros* (2da. edición). México, D. F.: McGraw-Hill.
- ▶▶ Walpole, R., Myers, R. y Myers, S. (2012). *Probabilidad y estadística para ingeniería y ciencias* (9na. edición). México, D. F.: Pearson Educación.