

# Asociación e independencia de variable categóricas

9/5/2020

## Contents

Distribución gamma y distribución exponencial . . . . .	1
Distribución de Ji cuadrada . . . . .	4
Aplicación de la prueba chi cuadrada . . . . .	4
Prueba de independencia . . . . .	5
Prueba de homogeneidad . . . . .	8

## Distribución gamma y distribución exponencial

Aunque la distribución normal se puede utilizar para resolver muchos problemas de ingeniería y ciencias, aún hay numerosas situaciones que requieren diferentes tipos de funciones de densidad. La distribución gamma deriva su nombre de la función gamma, que se define como:

Sea  $\Gamma : (0, \infty) \rightarrow \mathbf{R}$ , donde:

$$\Gamma(\alpha) = \int_{x=0}^{\infty} x^{\alpha-1} e^{-x} dx, \text{ para } \alpha > 0$$

Se integra por partes.  $u = x^{\alpha-1}$ ,  $dv = e^{-x}$  y se obtiene:

$$\Gamma(\alpha) = -e^{-x} x^{\alpha-1} \Big|_0^{\infty} + \int_{x=0}^{\infty} (\alpha-1) x^{\alpha-2} e^{-x} dx$$

Para  $\alpha > 1$ , lo que nos lleva a

$$\Gamma(\alpha) = (\alpha-1)\Gamma(\alpha-1)$$

Al aplicar reiteradamente la integración, se llega a:

$$\Gamma(\alpha) = (\alpha-1)(\alpha-2)\Gamma(\alpha-2) = (\alpha-1)(\alpha-2)(\alpha-3)\Gamma(\alpha-3)$$

, y así sucesivamente. Entonces, si  $\alpha = n$ , donde  $n$  es un entero positivo,

$$\Gamma(n) = (n-1)(n-2)\dots\Gamma(1)$$

Sin embargo, por definición  $\Gamma(1) = \int_{x=0}^{\infty} e^{-x} dx = 1$ , de donde

$$\Gamma(n) = (n-1)!$$

Otras propiedades de  $\Gamma(\alpha)$  son:

- \*  $\Gamma(n+1) = n!$  si  $n$  es un entero positivo
- \*  $\Gamma(n+1) = n\Gamma(n)$ ,  $n > 0$
- \*  $\Gamma(1/2) = \sqrt{\pi}$

La variable aleatoria continua  $X$  tiene una distribución gamma, con parámetros  $\alpha$  y  $\beta$ , si su función de densidad está dada por:

$$f(x : \alpha, \beta) = \begin{cases} \frac{1}{\beta^\alpha \Gamma(\alpha)} x^{\alpha-1} e^{-x/\beta} & \text{para } x > 0; \alpha, \beta > 0; \\ 0 & \text{en otro caso} \end{cases}$$

La distribución gamma especial para la que  $\alpha = 1$  se llama distribución exponencial.

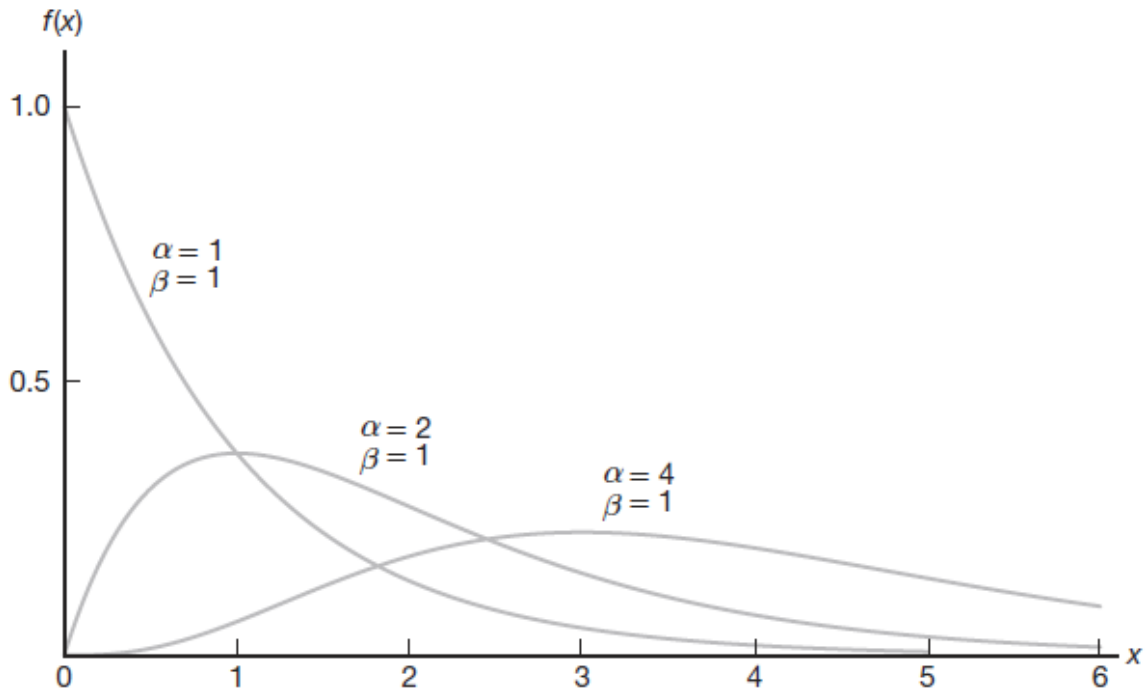


Figure 1: Distribuciones gamma de algunos valores de  $\alpha$  y  $\beta$ . La distribución gamma para  $\alpha = 1$  se llama distribución exponencial

La variable aleatoria continua  $X$  tiene distribución exponencial, con parámetro  $\beta$ , si su función de densidad está dada por:

$$f(x : \beta) = \begin{cases} \frac{1}{\beta} e^{-x/\beta} & \text{para } x > 0; \\ 0 & \text{en cualquier otro caso} \end{cases}$$

donde  $\beta > 0$

La media y la varianza de la distribución gamma son:

$$\mu = \alpha\beta \quad y \quad \sigma^2 = \alpha\beta^2$$

La media y la varianza de la distribución exponencial son:

$$\mu = \beta \quad y \quad \sigma^2 = \beta^2$$

Las aplicaciones más importantes de la distribución exponencial son situaciones donde se aplica el proceso de Poisson, que permite utilizar la distribución discreta llamada distribución de Poisson. La distribución de Poisson se utiliza para calcular la probabilidad de números específicos de eventos durante un período o espacio particulares. En muchas aplicaciones la variable aleatoria es el tiempo o la cantidad de espacio. La relación entre la distribución exponencial, también llamada exponencial negativa, y el proceso de Poisson es muy simple.

La distribución de Poisson tiene un parámetro  $\lambda$ , que se interpreta como el número medio de eventos por unidad de tiempo. Considere ahora la variable aleatoria descrita por el tiempo que se requiere para que ocurra el primer evento. Si se utiliza la distribución de Poisson, la probabilidad de que no ocurra algún evento, en el periodo hasta el tiempo  $t$ , es dada por:

$$p(0; \lambda t) = \frac{e^{-\lambda t} (\lambda t)^0}{0!} = e^{-\lambda t}$$

Si  $X$  es el tiempo para el primer evento de Poisson, entonces la probabilidad de que la duración del tiempo exceda a  $x$  es la misma probabilidad de que no ocurra algún evento de Poisson en  $x$ . Esto está dado por  $e^{-\lambda x}$ . Por lo tanto,

$$P(X > x) = e^{-\lambda x}$$

. De esta forma, la distribución acumulada para  $X$  está dada por:

$$P(0 \leq X \leq x) = 1 - e^{-\lambda x}$$

Si se deriva la función de distribución acumulativa anterior con el fin de obtener la función de densidad  $f(x) = \lambda e^{-\lambda x}$ , que es la función exponencial con  $\lambda = 1/\beta$ .

*Ejemplo de aplicación de distribución exponencial.*

Suponga que un sistema contiene cierto tipo de componente cuyo tiempo de operación antes de fallar, en años, está dado por  $T$ . La variable aleatoria  $T$  se modela bien mediante la distribución exponencial con tiempo medio de operación antes de fallar  $\beta = 5$ . Si se instalan 5 de estos componentes en diferentes sistemas, ¿cuál es la probabilidad de que al final de 8 años al menos dos aún funcionen?

La probabilidad de que un componente determinado siga funcionando después de 8 años es dada por:

$$P(T > 8) = \frac{1}{5} \int_8^{\infty} e^{-t/5} dt = e^{-8/5} \approx 0.2$$

Si  $X$  es el número de componentes que todavía funcionan después de 8 años. Entonces, utilizando la distribución binomial:

$$P(X \geq 2) = \sum_{x=2}^5 b(x; 5, 0.2) = 1 - \sum_{x=0}^1 b(x; 5, 0.2) = 1 - 0.7373 = 0.2627$$

```
x<- 2:5
sum(dbinom(x,5,0.2))
```

```
## [1] 0.26272
```

Suponga que las llamadas telefónicas que llegan a un conmutador particular siguen un proceso de Poisson con un promedio de 5 llamadas entrantes por minuto. ¿Cuál es la probabilidad de que transcurra hasta un minuto en el momento en que han entrado 2 llamadas al conmutador?

Se aplica el proceso de Poisson, con un lapso de tiempo hasta que ocurren 2 eventos de Poisson que sigue una distribución gamma con  $\beta = 1/5$  y  $\alpha = 2$ .

$$P(X \leq 1) = \int_0^1 1/\beta^2 x e^{-x/\beta} dx = 25 \int_0^1 x e^{-5x} dx = 1 - e^{-5}(1 + 5) \approx 0.92$$

```
x<-2
lambda<- 5
print(dist<-1-dpois(x,lambda))
```

```
## [1] 0.9157757
```

La distribución gamma trata con el tiempo hasta la ocurrencia de  $\alpha$  eventos de Poisson. También funciona aunque no exista una estructura de Poisson clara. Esto es particularmente aplicable para problemas de tiempo de supervivencia en aplicaciones de ingeniería y biomédicas.

*Ejemplo* En un estudio biomédico con ratas se utiliza una investigación de respuesta a la dosis para determinar el efecto de la dosis de un tóxico en su tiempo de supervivencia. Para cierta dosis del tóxico, el estudio determina que el tiempo de supervivencia de las ratas, en semanas, tiene una distribución gamma con  $\alpha = 5$  y  $\beta = 10$ . ¿Cuál es la probabilidad de que una rata no sobreviva más de 60 semanas?

$$P(X \leq 60) = \frac{1}{\beta^5} \int_0^6 \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(5)} dx$$

```
print(1-pgamma(60, 5, rate = 0.1, lower.tail = F))
```

```
## [1] 0.7149435
```

## Distribución de Ji cuadrada

Un caso especial de la distribución gamma se obtiene al permitir que  $\alpha = \nu/2$  y  $\beta = 2$ , donde  $\nu$  es un entero positivo. Este resultado se conoce como distribución chi cuadrada. La distribución tiene un solo parámetro,  $\nu$ , denominado grados de libertad.

$$f(x : \nu) = \begin{cases} \frac{1}{2^{\nu/2} \Gamma(\nu/2)} x^{\nu/2-1} e^{-x/2} & \text{para } x > 0 \\ 0 & \text{en otro caso} \end{cases}$$

donde  $\nu$  es positiva.

La media y la varianza de la distribución  $\chi^2$  son

$$\mu = \nu \quad y \quad \sigma^2 = 2\nu$$

## Aplicación de la prueba chi cuadrada

La prueba chi-cuadrada, también llamada Ji cuadrado ( $\chi^2$ ), se encuentra dentro de las pruebas pertenecientes a la estadística descriptiva aplicada al estudio de dos variables.

*Nota:* La estadística descriptiva se centra en extraer información sobre la muestra, la estadística inferencial extrae información sobre la población.

La distribución chi cuadrada no es simétrica y depende de un número específico de grados de libertad (gl). Los gl adecuados para el estadístico  $\chi^2$  dependen de la aplicación para la que se utilice.

### Grados de libertad

Los grados de libertad (gl) de un estadístico calculado con base en  $n$  datos, se refiere al número de cantidades independientes que se necesitan en su cálculo, menos el número de restricciones que relacionan a las observaciones y el estadístico.

## Prueba de probabilidades de valores esperados

Los valores esperados para cada categoría se calculan usando probabilidades hipotéticas,  $E_i = np_i$ , se usan para calcular el valor observado del estadístico de prueba  $\chi^2$ . Para un experimento multinomial formado por  $k$  categorías o celdas, el estadístico de prueba tiene una distribución  $\chi^2$  aproximada, con  $gl = (k - 1)$ .

### Ejemplo

Un investigador diseña en el que una rata es atraída al final de una rampa que se divide, llevando a tres

puertas de colores diferentes. Se hace que una rata baje 90 veces y se observan las siguientes elecciones. ¿La rata tiene preferencia por alguna de las tres puertas?

Cantidad Observada ( $O_i$ )

Puerta verde	Puerta roja	Puerta azul
20	39	31

$H_0 : p_1 = p_2 = p_3 = 1/3$   $H_a$  = Al menos una  $p_i$  es diferente de  $1/3$

$p_i$  es la probabilidad de que la rata elija la puerta  $i$ , para  $i = 1, 2$  y  $3$ . Los valores esperados son los mismos para cada una de las categorías  $np_i = 90(1/3) = 30$

El estadístico de prueba  $\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i} = (20 - 30)^2/30 + (39 - 30)^2/30 + (31 - 30)^2/30 = 6.067$

En este caso se tienen  $n - 1 = 2$  grados de libertad pues la única restricción lineal en las probabilidades es que su suma es igual a 1.

El valor observado para  $\chi^2 = 6.067$  está entre 5.99 y 7.39 en la tabl, el valor de  $p$  está entre 0.025 y 0.05. EL valor que se toma como referencia en investigación es, generalmente,  $p < 0.05$ .

Esto significa que la hipótesis nula de no preferencia por alguna puerta se rechaza. La rata prefiere alguna de las puertas.

Puerta azul  $31/90 = 0.344$

Puerta verde  $20/90 = 0.222$

Puerta roja  $39/90 = 0.433$

La rata tiene preferencia por la puerta roja, pero no se puede declarar una relación **causal** entre el color y la preferencia.

## Prueba de independencia

La prueba chi-cuadrado es una de las más conocidas y utilizadas para analizar variables nominales o cualitativas, es decir, para determinar la existencia o no de independencia entre dos variables. Que dos variables sean independientes significa que no tienen relación, y que por lo tanto una no depende de la otra, ni viceversa.

*Ejemplo* Suponga que se desea determinar si las opiniones de los votantes residentes del estado de Illinois respecto a una nueva reforma fiscal son independientes de sus niveles de ingreso. Los sujetos de una muestra aleatoria de 1000 votantes registrados del estado de Illinois se clasifican de acuerdo con su posición en las categorías de ingreso bajo, medio o alto, y si están a favor o no de la nueva reforma fiscal. Las frecuencias observadas se presentan en la siguiente tabla, que se conoce como tabla de contingencia.

Reforma Fiscal	Nivel bajo	Nivel medio	Nivel Alto
A favor	182	213	203
En contra	154	138	110
Total	336	351	313

Esta es una tabla de  $2 \times 3$

Una tabla de contingencia con  $r$  renglones y  $c$  columnas se denomina tabla  $r \times c$  (" $r \times c$ " se lee " $r$  por  $c$ "). Los totales de renglones y columnas en la tabla se denominan frecuencias marginales.

La decisión de aceptar o rechazar la hipótesis nula,  $H_0$ , de que la opinión de un votante respecto a la nueva reforma fiscal es independiente de su nivel de ingreso, se basa en qué tan bien se ajusten las frecuencias

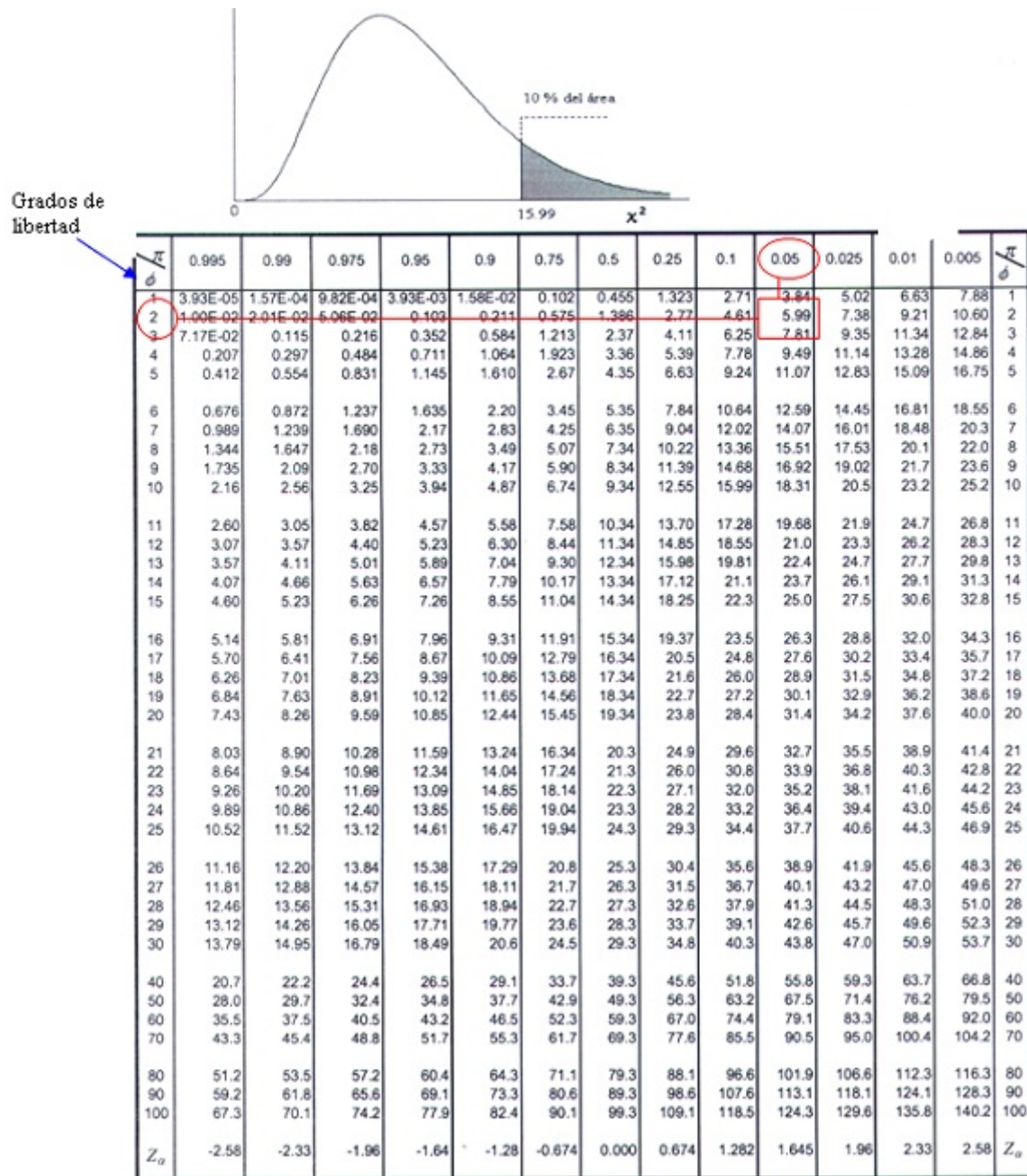


Figure 2: Tabla de valores críticos de  $\chi^2$

observadas en cada una de las 6 celdas de la tabla de contingencia y en las frecuencias esperadas si supusiéramos que  $H_0$  es verdadera. Para encontrar estas frecuencias esperadas se definen los siguientes eventos:

L: Una persona seleccionada está en el nivel de ingresos bajo.

M: Una persona seleccionada está en el nivel de ingresos medio.

H: Una persona seleccionada está en el nivel de ingresos alto.

F: Una persona seleccionada está a favor de la nueva reforma fiscal.

A: Una persona seleccionada está en contra de la nueva reforma fiscal.

Con base en las frecuencias marginales:

$$P(L) = 336/1000$$

$$P(M) = 351/1000$$

$$P(H) = 313/1000$$

$$P(F) = 598/1000$$

$$P(A) = 402/1000$$

Si  $H_0$  es verdadera y las dos variables son independientes:

$$P(L \cap F) = P(L)P(F) = (336/1000)(598/1000)$$

$$P(L \cap A) = P(L)P(A) = (336/1000)(402/1000)$$

$$P(M \cap F) = P(M)P(F) = (351/1000)(598/1000)$$

$$P(M \cap A) = P(M)P(A) = (351/1000)(402/1000)$$

$$P(H \cap F) = P(H)P(F) = (313/1000)(598/1000)$$

$$P(H \cap A) = P(H)P(A) = (313/1000)(402/1000)$$

Las frecuencias esperadas se obtienen multiplicando la probabilidad de cada celda por el número total de observaciones. De esta manera, se estima que el número esperado de votantes de bajo ingreso en la muestra que favorecen la reforma fiscal cuando  $H_0$  es verdadera:

$$(336/1000)(598/1000)(1000) = 200.9$$

Esto es,

$$frecuenciaesperada = \frac{(total \quad columna)(total \quad renglón)}{gran \quad total}$$

### Actividad:

completa la tabla con la frecuencia esperada para cada “celda”:

Reforma Fiscal	Nivel bajo	Nivel medio	Nivel Alto	Total
A favor	182 (200.9)	213()	203()	598
En contra	154()	138()	110()	402
Total	336	351	313	1000

El número de grados de libertad asociados con la prueba chi cuadrada que aquí se usa es igual al número de frecuencias de celdas que se pueden llenar libremente cuando se proporcionan los totales marginales y el gran total, y en este caso es 2.

Una fórmula sencilla que proporciona el número correcto de grados de libertad es  $\nu = (r - 1)(n - 1) =$

$(2 - 1)(3 - 1) = 2$ . Calcule

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

donde la sumatoria se extiende a todas las celdas  $rc$  en la tabla de contingencia  $r \times c$

Si  $\chi^2 > \chi_\alpha$  con  $n$  grados de libertad, rechace la hipótesis nula de independencia al nivel de significancia  $\alpha$ , en otro caso acéptela.

$$\chi^2 = 7.85$$

con  $P \approx 0.02$  *Actividad: demuestre este resultado*

En la tabla de valores para  $\chi^2$  se encuentra que  $\chi_{0.05}^2$  para 2 gl.

**Rechazamos la hipótesis nula y concluimos que la opinión de un votante respecto a la reforma fiscal y su nivel de ingresos NO son independientes.** Los valores  $\chi^2$  calculados dependen de las frecuencias de las celdas y, en consecuencia, son discretos.

La distribución chi cuadrada continua parece aproximarse muy bien a la distribución de muestreo discreta de  $\chi^2$ , siempre y cuando el número de grados de libertad sea mayor que 1.

### Corrección de Yates

En una tabla de contingencia de  $2 \times 2$ , donde sólo se tiene 1 grado de libertad, se aplica una corrección llamada corrección de Yates para continuidad. La fórmula corregida entonces se convierte en

$$\chi^2(\text{corregida}) = \sum_i \frac{(|o_i - e_i| - 0.5)^2}{e - i}$$

Si las frecuencias de las celdas esperadas son grandes, los resultados corregidos y sin corrección son casi iguales. Cuando las frecuencias esperadas están entre 5 y 10, se debe aplicar la corrección de Yates.

### Prueba de homogeneidad

Otro tipo de problema para el que se aplica el método anterior es aquel en el cual los totales de renglón y de columna están predeterminados.

#### Ejemplo

Suponga que decidimos de antemano seleccionar 200 demócratas, 150 republicanos y 150 independientes entre los votantes del estado de Carolina del Norte y registrar si están a favor de una iniciativa de ley para el aborto, si están en contra o si están indecisos. |Ley para el aborto |Demócrata |Republicano |Independiente| Total|  
| ——— |—————| ———|—————| :————: |A favor|En contra|Indeciso|Total 82|93|25|214 70|62|18|222 62|67|21|64  
200|150|150|500 Ahora bien, en vez de hacer una prueba de independencia, probamos la hipótesis de que las proporciones de población dentro de cada renglón son iguales. Es decir, probamos la hipótesis de que las proporciones de demócratas, republicanos e independientes que están a favor de la ley para el aborto son iguales; las proporciones de cada afiliación política contra la ley son iguales y las proporciones de cada afiliación política que están indecisos son iguales.

Se busca determinar si las tres categorías de votantes son homogéneas en lo que se refiere a sus opiniones acerca de la iniciativa de ley para el aborto. A esta prueba se le conoce como prueba de homogeneidad.

Al suponer homogeneidad de nuevo calculamos las frecuencias esperadas de las celdas multiplicando los totales de renglón y de columna correspondientes y después dividiendo entre el gran total. Luego continuamos el análisis utilizando el mismo estadístico chi cuadrada como antes.

1.  $H_0$ : Para cada opinión las proporciones de demócratas, republicanos e independientes son iguales.
2.  $H_1$ : Para al menos una opinión las proporciones de demócratas, republicanos e independientes no son iguales.
3.  $\alpha = 0.05$ .
4. Región crítica:  $\chi^2 > 9.488$  con  $v = 4$  grados de libertad (de la tabla de  $\chi^2$ ).



5. Calcular las 4 frecuencias de las celdas usando la fórmula de las frecuencias de las celdas esperadas  $frecuenciaesperada = \frac{(total\ columna)(total\ renglón)}{gran\ total}$ . Todas las demás frecuencias se obtienen mediante sustracción.

$\chi^2 = 1.53$  *Actividad: demuestre este resultado* Conclusión: No rechazar  $H_0$ . No hay suficiente evidencia para concluir que la proporción de demócratas, republicanos e independientes difiere para cada opinión expresada.

*Actividades*

1. Individual: Resolver los ejercicios de práctica en la plataforma Connect “Aplicación chi cuadrada”.
2. En grupos de tres alumnos: Para la base de datos BD\_polimorf.xlsx, que está en Moodle, determina la independencia de las variables HTA, ALCOHOL y TABQ con respecto a POL1, POL2 y POL3. En tres ejercicios separados.