

# Análisis de Datos Categóricos

Leticia Gracia Medrano



# Contents

<b>1</b>	<b>Introducción</b>	<b>1</b>
1.1	Ejemplos de donde salen datos categóricos . . . . .	1
1.2	Modelo Poisson . . . . .	1
1.3	Modelo Binomial . . . . .	2
1.4	Modelo Multinomial . . . . .	3
1.5	Inferencia sobre $\pi$ . . . . .	3
<b>2</b>	<b>Tablas de contingencia</b>	<b>7</b>
2.1	Notación . . . . .	7
2.2	Esquemas de muestreo . . . . .	8
2.2.1	Esquema Poisson . . . . .	8
2.2.2	Esquema Multinomial . . . . .	8
2.2.3	Esquema Multinomial-Producto . . . . .	9
2.2.4	Hipótesis de no asociación . . . . .	10
2.3	Prueba de Independencia . . . . .	10
2.4	La $\chi^2$ . . . . .	12
2.5	Hago una tabla de 2 por 2 en R . . . . .	13
2.5.1	Ejemplo 1 . . . . .	13
2.5.2	Ejemplo 2 . . . . .	14
2.5.3	Ejemplo 3 . . . . .	15
2.6	Más de la Ji Cuadrada . . . . .	16

2.6.1	Sensible al tamaño de muestra . . . . .	16
2.7	Comparación de proporciones en tablas de 2 por 2 . . . . .	18
2.8	Diferencia de proporciones . . . . .	18
2.9	Riesgo Relativo . . . . .	18
2.10	Cociente de Momios . . . . .	19
2.11	Residuales . . . . .	22
2.12	Ejemplo . . . . .	22
2.13	Intervalo de confianza para $\log(\theta)$ . . . . .	25
2.14	Intervalo de confianza para $\log(RR)$ . . . . .	27
2.14.1	Ejemplo caso y controles . . . . .	27

# Chapter 1

## Introducción

### 1.1 Ejemplos de donde salen datos categóricos

En muchas ocasiones analizamos variables como: nacionalidad, tipo de escuela a la que asiste, marca de su preferencia, partido por el que se votaría, que son **nominales** y otras parecidas como: escolaridad, calificaciones de algún servicio, que son **ordinales**. Estas son las variables que analizaremos en este curso. En el primer caso se pueden intercambiar las categorías de las variables y no se pierde información, mientras que en el segundo caso debemos mantener las etiquetas pues éstas si mantienen un **orden**. También conviene hacer diferencia entre variable *respuesta o dependiente* ( $Y$ ) y las variables *explicativas o independientes* ( $X$ ). En este curso el interés está puesto en las variables respuesta categóricas, las explicativas podrán ser continuas o categóricas, según sea el contexto de estudio.

### 1.2 Modelo Poisson

La distribución Poisson ayuda a modelar conteos como: número de accidentes en un tramo carretero, personas que llegan a formarse en una fila. La distribu-

ción está dada por

$$P(y) = \frac{\exp^{-\mu} \mu^y}{y!}$$

Por ejemplo: Si se tiene que ocurren un promedio de 2 accidentes semanales en cierto tramo carretero entonces la probabilidad de tener 0 accidentes en una semana dada es:

$$P(0) = \frac{\exp^{-2} 2^0}{0!} = \exp^{-2} = 0.135$$

En esta distribución  $E(Y) = Var(Y) = \mu$ , aquí entonces si la media aumenta, también aumenta la varianza. En la práctica ocurre frecuentemente que los conteos tienen mayor varianza que la esperada, esto se conoce como *sobredispersión*. En ocasiones suponer la distribución Poisson resulta muy simplista, pero en otras resulta muy útil.

### 1.3 Modelo Binomial

En el ejemplo anterior el número de accidentes es aleatorio. Pero podría plantearse algo como que se clasifican los accidentes hasta que ocurren  $N$ , con el propósito de estimar la proporción de estos que resultan fatales, entonces el total de accidentes es fijo. Ahora el número de accidentes fatales ya no es Poisson porque tiene un tope máximo de  $N$ .

Si se tienen que el número de accidentes fatales en  $t$  semanas tiene una media de  $2t$ , y la tasa para accidentes no-fatales es de  $8t$ . Cuando se junta un total de  $N$  accidentes ocurre que el número de accidentes fatales se distribuye como binomial con parámetros  $N$  y  $\pi = \frac{2t}{2t+8t} = .2$ , la probabilidad de cualquier accidente resulte fatal. La función de distribución binomial recordemos está dada por:

$$P(y) = \frac{N!}{y!(N-y)!} \pi^y (1-\pi)^{N-y} \quad \text{con } y = 0, 1, 2, \dots, N$$

Para el caso en que  $N = 10$  y  $\pi = .2$  la probabilidad de que haya  $y = 0$  accidentes es

$$P(0) = \frac{10!}{0!(10)!} \cdot 2^0 (.8)^{10} = (.8)^{10} = .107;$$

Para esta distribución  $E(Y) = N\pi$  y  $Var(Y) = N\pi(1 - \pi)$ , entonces aquí la varianza siempre es menor que la media.

Cuando los resultados pueden ser más que dos, se tiene una distribución *multinomial*, que se verá en el último tema.

## 1.4 Modelo Multinomial

En este caso el experimento tiene  $c$  posibles resultados, sus probabilidades las denotamos por  $\pi_1, \pi_2, \dots, \pi_c$ , donde  $\sum \pi_j = 1$ .

Para  $n$  observaciones independientes, la probabilidad multinomial de que  $n_1$  caigan en la categoría 1,  $n_2$  caigan en la categoría 2, ... y  $n_c$  caigan en la categoría  $c$ , donde  $\sum n_j = n$  es

$$P(n_1, n_2, \dots, n_c) = \left( \frac{n!}{n_1! n_2! \dots n_c!} \right) \pi_1^{n_1} \pi_2^{n_2} \dots \pi_c^{n_c}$$

La distribución multinomial es multivariada. La marginal para cualquier categoría es binomial. Para la categoría  $j$  el conteo  $n_j$  tiene media  $n\pi_j$  y desviación estándar  $\sqrt{n\pi_j(1 - \pi_j)}$

## 1.5 Inferencia sobre $\pi$

En el modelo binomial el parámetro es  $\pi$  que generalmente es desconocido y a través de una muestra trataremos de estimarlo.

La probabilidad de los datos observados, expresada como una función del parámetro, es la función de verosimilitud. Para una  $n = 10$  y una  $y = 0$  la función de verosimilitud es:

$$P(0) = (10!/0!10!) \pi^0 (1 - \pi)^{10} = (1 - \pi)^{10}$$

$$\ell(\pi) = (1 - \pi)^{10}$$

que alcanza el máximo cuando  $\pi = 0$ , entonces el resultado  $y = 0$  ocurre con una mayor probabilidad cuando  $\pi = 0$ . OJO esta función depende de los valores que tome  $\pi$ , para cada valor de  $y$  se tiene una función de verosimilitud distinta, así para  $y = 6$

$$\ell(\pi) = (10!/6!4!)\pi^6(1 - \pi)^4 = 210 * \pi^6(1 - \pi)^4$$

cuyo máximo lo alcanza en  $\pi = .6$ , en este caso  $y = 6$  ocurre con una mayor probabilidad cuando  $\pi = .6$

Se tiene que el estimador máximo verosímil para  $\pi$  es  $\hat{p} = \sum_{i=1}^n x_i/n$  donde  $x_i$  es 1 o 0 según se observe éxito o fracaso. Entonces  $\hat{p}$  es un promedio por lo que se puede utilizar el Teorema Central del Límite. Para una  $n$  grande  $y = \bar{x}$  se distribuye aproximadamente como una normal con media  $E(\hat{p}) = \sum E(x_i)/n = n\pi/n = \pi$  y varianza  $var(\hat{p}) = \frac{\sum var(x_i)}{n^2} = \frac{n\pi(1-\pi)}{n^2} = \pi(1 - \pi)/n$ .

Se puede usar la estadística

$$z = \frac{\hat{p} - \pi_o}{\sqrt{\pi_o(1 - \pi_o)/n}}$$

para probar la hipótesis nula  $H_o : \pi = \pi_o$ .

Y se puede construir un intervalo de  $100(1 - \alpha)$  de confianza para  $\pi$  con

$$\hat{p} \pm z_{\alpha/2}SE, \text{ con } SE = \sqrt{\hat{p}(1 - \hat{p})/n}$$

donde  $z_{\alpha/2}$  es el percentil que deja a la derecha una cola tamaño  $\alpha/2$ .

Hay que tener cuidado, esta aproximación es buena cuando  $\pi$  está cerca de 0.5 o cuando la  $n$  es muy grande. Si no se tiene eso, el nivel de confianza disminuiría. Y es muy mala cuando  $\pi$  se acerca a uno o al cero.

Que más se podría hacer?? Si existe una correspondencia entre pruebas de hipótesis e intervalos de confianza, que dice que el intervalo de confianza es



aquel donde la hipótesis nula No se rechaza, o sea la región no crítica. Entonces deben hallarse los valores de  $\pi_o$  de manera que:

$$\frac{|\hat{p} - \pi_o|}{\sqrt{(\pi_o(1 - \pi_o)/n)}} = 1.96 = z_{.05/2}$$

Cooooo??? pues elevando al cuadrado ambos lados de la ecuación anterior y de la que resulta una ecuación cuadrática en  $\pi_o$  Para el caso de  $\hat{p} = 0.90$  y  $n = 10$  usando esto se encuentra las raíces  $\pi_{o1} = .596$  y  $\pi_{o2} = .982$ , que darían un intervalo  $(.596, .982)$ , mientras que usando la ecuación "tradicional" se tendría un intervalo al 90% de confianza  $(0.714, 1.086)$ .

Otra aproximación es la llamada de Agresti-Coull que es muy fácil, se suman 2 a los éxitos y se suman 2 a los fracasos y entonces ... usando el ejemplo anterior

$\hat{p} = (9 + 2)/(10 + 2 + 2) = .786$  y  $SE = (.786)(.214)/14 = .110$  con lo que se obtiene un intervalo de  $(.57, 1.0)$ .

Este método funciona bien aún con muestras pequeñas.



# Chapter 2

## Tablas de contingencia

### 2.1 Notación

Cuando se tienen dos criterios o variables de clasificación de las observaciones, al hacer el “cruce” de éstas se genera una tabla de frecuencias como sigue:

$$\begin{bmatrix} n_{11} & n_{12} & \cdots & n_{1j} & \cdots & n_{1J} \\ n_{21} & n_{22} & \cdots & n_{2j} & \cdots & n_{2J} \\ \vdots & \vdots & & \vdots & & \vdots \\ n_{i1} & n_{i2} & \cdots & n_{ij} & \cdots & n_{iJ} \\ \vdots & \vdots & & \vdots & & \vdots \\ n_{I1} & n_{I2} & \cdots & n_{Ij} & \cdots & n_{IJ} \end{bmatrix}$$

Donde la variable renglón tiene  $I$  categorías y la variables columna tiene  $J$  categorías.

Los totales por columna son:

$$\sum_{i=1}^I n_{ij} = n_{.j}$$

los totales por renglón:

$$\sum_{j=1}^J n_{ij} = n_i.$$

el total general :

$$\sum_{j=1}^J \sum_{i=1}^I n_{ij} = n$$

y las frecuencias relativas:  $r_i = n_{i.}/n$  y  $c_j = n_{.j}/n$

Con  $i = 1, \dots, I$  y  $j = 1, \dots, J$

## 2.2 Esquemas de muestreo

En el análisis de tablas de contingencia de dos dimensiones se utilizan tres esquemas de muestreo que ocurren en la práctica: Esquema Poisson, Esquema Multinomial, y Esquema Multinomial-Producto.

### 2.2.1 Esquema Poisson

Supone que se observa cualquier cantidad de datos  $\{n_{ij}\}$  durante un intervalo de tiempo. La distribución marginal de las observaciones  $n_{ij}$  es  $Poisson(m_{ij})$ , para  $i = 1, \dots, I$ ,  $j = 1, \dots, J$ . Las probabilidades marginales son

$$P(n_{ij}|m_{ij}) = e^{-m_{ij}} \frac{m_{ij}^{n_{ij}}}{n_{ij}!},$$

donde  $\sum_{i=1}^I \sum_{j=1}^J n_{ij} = n_{..}$ ,  $E(n_{ij}) = m_{ij}$ .

La función de verosimilitud queda expresada como

$$L(m) = L(\{m_{ij}\}) = \prod_{i=1}^I \prod_{j=1}^J e^{-m_{ij}} \frac{m_{ij}^{n_{ij}}}{n_{ij}!}.$$

### 2.2.2 Esquema Multinomial

Este esquema supone que el número total de observaciones  $n_{..} = \sum_{i=1}^I \sum_{j=1}^J n_{ij}$  es fijo. Se tienen  $I \times J$  categorías con probabilidad  $Multinomial(n_{..}, \{p_{ij}\})$ , en donde las probabilidades están dadas por

$$P(n_{ij}|m_{ij}) = \binom{n_{..}}{n_{11}n_{12}\dots n_{IJ}} \prod_{i=1}^I \prod_{j=1}^J p_{ij}^{n_{ij}} = \frac{n_{..}!}{\prod_{i=1}^I \prod_{j=1}^J n_{ij}!} \prod_{i=1}^I \prod_{j=1}^J p_{ij}^{n_{ij}},$$

donde  $\sum_{i=1}^I \sum_{j=1}^J n_{ij} = n_{..}$ ,  $\prod_{i=1}^I \prod_{j=1}^J p_{ij} = 1$ ,  $m_{ij} = E(n_{ij}) = n_{..} p_{ij}$ , y  $p_{ij} = \frac{m_{ij}}{n_{..}}$ . Además  $m_{..} = \sum_{i=1}^I \sum_{j=1}^J m_{ij} = \sum_{i=1}^I \sum_{j=1}^J n_{..} p_{ij} = n_{..}$ .

La función verosimilitud es

$$\begin{aligned}
 L(m) &= L(\{m_{ij}\}) \\
 &= \frac{n_{..}!}{\prod_{i=1}^I \prod_{j=1}^J n_{ij}!} \prod_{i=1}^I \prod_{j=1}^J p_{ij}^{n_{ij}} \quad \text{dado que } p_{ij} = \frac{m_{ij}}{n_{..}} \\
 &= \frac{n_{..}!}{\prod_{ij} n_{ij}!} \prod_{ij} \left( \frac{m_{ij}}{n_{..}} \right)^{n_{ij}} \\
 &= \frac{n_{..}!}{\prod_{ij} n_{ij}!} \prod_{ij} \left( \frac{m_{ij}}{n_{..}} \right)^{n_{ij}} \frac{e^{-\sum_{ij} m_{ij}}}{e^{-m_{..}}} \quad \text{dado que } n_{..} = m_{..} \\
 &= \frac{n_{..}!}{e^{-n_{..}} \prod_{ij} n_{..}^{n_{ij}}} \prod_{ij} \left( e^{-m_{ij}} \frac{m_{ij}^{n_{ij}}}{n_{ij}!} \right) \\
 &\propto \prod_{ij} e^{-m_{ij}} \frac{m_{ij}^{n_{ij}}}{n_{ij}!},
 \end{aligned}$$

que indica que el esquema Multinomial es equivalente al esquema Poisson.

### 2.2.3 Esquema Multinomial-Producto

En este esquema se supone que los totales marginales por renglón (o columna)  $n_{1.}, n_{2.}, \dots, n_{I.}$  (o  $n_{.1}, n_{.2}, \dots, n_{.J}$ ) están fijos. Para estos, en cada renglón se tiene una distribución multinomial. Entonces, para el renglón  $i$ ,  $i = 1, \dots, I$ , las probabilidades son

$$P(n_{ij}|m_{ij}) = \binom{n_{i.}}{n_{i1} n_{i2} \dots n_{iJ}} \prod_{j=1}^J p_{j|i}^{n_{ij}} = \frac{n_{i.}!}{\prod_{j=1}^J n_{ij}!} \prod_{j=1}^J p_{j|i}^{n_{ij}},$$

donde  $\sum_{i=1}^I n_{i.} = n_{..}$ ,  $\prod_{j=1}^J p_{j|i} = 1$ ,  $m_{ij} = E(n_{ij}) = n_{i.} p_{j|i}$ , y  $p_{j|i} = \frac{m_{ij}}{n_{i.}}$ .

Además. La verosimilitud es

$$\begin{aligned}
 L(m) &= L(\{m_{ij}\}) \\
 &= \prod_{i=1}^I \left[ \frac{n_{i\cdot}!}{\prod_{j=1}^J n_{ij}!} \prod_{j=1}^J p_{j|i}^{n_{ij}} \right] \quad \text{dado que } p_{j|i} = \frac{m_{ij}}{n_{i\cdot}} \\
 &= \frac{\prod_i n_{i\cdot}!}{\prod_{ij} n_{ij}!} \prod_{ij} \left( \frac{m_{ij}}{n_{i\cdot}} \right)^{n_{ij}} \\
 &= \frac{\prod_i n_{i\cdot}!}{\prod_{ij} n_{ij}!} \prod_{ij} \left( \frac{m_{ij}}{n_{i\cdot}} \right)^{n_{ij}} \frac{e^{-\sum_{ij} m_{ij}}}{e^{-m_{\cdot\cdot}}} \quad \text{dado que } n_{\cdot\cdot} = m_{\cdot\cdot} \\
 &= \frac{\prod_i n_{i\cdot}!}{e^{-n_{\cdot\cdot}} \prod_{ij} n_{i\cdot}^{n_{ij}}} \prod_{ij} \left( e^{-m_{ij}} \frac{m_{ij}^{n_{ij}}}{n_{ij}!} \right) \\
 &\propto \prod_{ij} e^{-m_{ij}} \frac{m_{ij}^{n_{ij}}}{n_{ij}!},
 \end{aligned}$$

que indica que el esquema Multinomial-Producto es equivalente al esquema Poisson.

### 2.2.4 Hipótesis de no asociación

La hipótesis nula en la Ji-Cuadrada ( $\chi^2$ ) de no asociación corresponde a “diferentes interpretaciones” para cada esquema de muestreo:

- Esquema Poisson: No Asociación (las variables no están relacionadas).
- Esquema Multinomial: Independencia (la probabilidad conjunta es el producto de las probabilidades marginales).
- Esquema Multinomial-Producto: Homogeneidad (la distribución es la misma en cada renglón).

## 2.3 Prueba de Independencia

La pregunta que surge es si las dos variables son **independientes**, es decir si los datos se acomodan en la tabla de manera proporcional al total de los

renglones y el total de las columnas. Si los datos no se acomodan de manera proporcional diremos que ciertas categorías de las variables están **asociadas**. Es claro que las proporciones no resultan exactas sino que hay variaciones aleatorias, pero si las diferencias son muy grandes con esas proporciones “esperadas”, se dirá que las variables no son independientes.

Si en la población de la que se saca la muestra, la probabilidad de que una observación pertenezca a la celda  $i, j$  se llama  $p_{i,j}$ , entonces la frecuencia esperada  $F_{i,j}$  de observaciones para esa celda, luego de sacar una muestra tamaño  $N$  es  $F_{i,j} = n \times p_{i,j}$ .

Ahora, si  $p_{i.}$  es la probabilidad de pertenecer al renglón  $i$  y  $p_{.j}$  es la probabilidad de pertenecer a la columna  $j$ , cuando las variables son independientes ocurre que  $p_{ij} = p_{i.} \times p_{.j}$

Entonces las frecuencias esperadas cuando las variables son independientes son:

$$F_{ij} = N \times p_{i.} \times p_{.j}$$

Estas probabilidades no se conocen, pero pueden ser estimadas con

$$\hat{p}_{i.} = \frac{n_{i.}}{n}$$

y

$$\hat{p}_{.j} = \frac{n_{.j}}{n}$$

y las frecuencias esperadas se estiman con

$$E_{ij} = n\hat{p}_{i.}\hat{p}_{.j} = n \frac{n_{i.}}{n} \frac{n_{.j}}{n} = \frac{n_{i.}n_{.j}}{n}$$

Hay ocasiones que los totales por renglón son fijos, esto por el diseño del muestreo, si la variables respuesta  $Y$  es binaria se tiene un modelo binomial, si tiene más categorías se tiene un esquema multinomial. Y en ese caso nos fijamos en las distribuciones condicionales para cada nivel dela variable  $X$ . Aquí la independencia entre  $X$  y  $Y$  puede expresarse también como que las

distribuciones condicionales de  $Y$  para cada nivel de la variable  $X$  son las mismas.

Otra situación es cuando  $n$  es fija y clasificamos a los individuos al “cruzar” las dos variables respuesta, en ese caso se tiene una distribución multinomial con  $I \times J$  categorías.

## 2.4 La $\chi^2$

Para analizar si las variables son independientes se puede utilizar la estadística  $\chi^2$ , dada por

$$\chi^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(n_{ij} - \frac{n_{i.}n_{.j}}{n})^2}{\frac{n_{i.}n_{.j}}{n}} = \sum_{i=1}^I \sum_{j=1}^J \frac{(n_{ij} - E_{ij})^2}{E_{ij}}$$

Cuando esta estadística toma valores “grandes” se rechaza la hipótesis nula de independencia.

La distribución asintótica de esta estadística puede hallarse suponiendo que las frecuencias observadas siguen una distribución multinomial y que las frecuencias esperadas no son muy pequeñas, y corresponde a una distribución **ji cuadrada** con  $(I - 1) * (J - 1)$  grados de libertad.

Muchos paquetes trabajan con la Ji Cuadrada con corrección de Yates, esto es:

$$\chi^2 = \sum_{ij} [ |o_{ij} - e_{ij}| - 0.5 ]^2 / e_{ij}$$



## 2.5 Hago una tabla de 2 por 2 en R

### 2.5.1 Ejemplo 1

Haciendo esto muy rápido nada más doy los datos (pero no se ve bonito).

```
> v1 <- matrix(c(25, 11, 12, 14), ncol = 2)
```

```
> chisq.test(v1)
```

```
      Pearsons Chi-squared test with Yates continuity correction
```

```
data:  v1
```

```
X-squared = 2.5041, df = 1, p-value = 0.1135
```

En este caso para un  $\alpha = .10$  la prueba resulta no significativa, es decir no rechazamos  $H_0$ , es decir no hay evidencia para decir que no son independientes

Si desean hacer una salida más bonita ponen:

```
> FUMA <- c("NO FUMA", "NO FUMA", "SI FUMA", "SI FUMA")
```

```
> GENERO <- c("FEM", "MASC", "FEM", "MASC")
```

```
> conteos <- c(25, 11, 12, 14)
```

```
> TABLA <- data.frame(FUMA, GENERO, conteos)
```

```
> xtabs(conteos ~ FUMA + GENERO, data = TABLA)
```

```

      GENERO
FUMA    FEM MASC
NO FUMA  25  11
SI FUMA  12  14
```

```
> chisq.test(xtabs(conteos ~ FUMA + GENERO, data = TABLA))
```

Pearsons Chi-squared test with Yates continuity correction

```
data: xtabs(conteos ~ FUMA + GENERO, data = TABLA)
X-squared = 2.5041, df = 1, p-value = 0.1135
```

## 2.5.2 Ejemplo 2

La proporción de niños de bajo peso al nacer es la misma en las mujeres que fuman, que en las que no fuman?

```
> FUMA <- c("NO FUMA", "NO FUMA", "SI FUMA", "SI FUMA")
> BAJOPESO <- c("SI", "NO", "SI", "NO")
> conteos <- c(105, 1645, 43, 207)
> TABLA <- data.frame(FUMA, BAJOPESO, conteos)
> xtabs(conteos ~ FUMA + BAJOPESO, data = TABLA)
```

	BAJOPESO	
FUMA	NO	SI
NO FUMA	1645	105
SI FUMA	207	43

```
> chisq.test(xtabs(conteos ~ FUMA + BAJOPESO, data = TABLA))
```

Pearsons Chi-squared test with Yates continuity correction

```
data: xtabs(conteos ~ FUMA + BAJOPESO, data = TABLA)
X-squared = 38.4266, df = 1, p-value = 5.685e-10
```

La prueba es altamente significativa. Decimos que la proporción de niños de bajo peso es estadísticamente diferente en el grupo de las mamás que fuman ( $43/250=0.172$ ) que en las que no fuman ( $105/1645=0.0638$ ) . .

### 2.5.3 Ejemplo 3

La hipótesis que se desea probar con los datos de la siguiente tabla, es si el tipo de tuberculosis por el que la persona muere es independiente del género

```
> tipotuberculosis <- c("resp", "otra", "resp", "otra")
> genero <- c("m", "m", "f", "f")
> conteos <- c(3534, 270, 1319, 252)
> TABLA <- data.frame(tipotuberculosis, genero, conteos)
> xtabs(conteos ~ tipotuberculosis + genero, data = TABLA)
```

```
          genero
tipotuberculosis  f    m
          otra  252  270
          resp 1319 3534
```

```
> tabla3 <- xtabs(conteos ~ tipotuberculosis + genero, data = TABLA)
```

Para calcular los marginales por rengl\on

```
> margin.table(tabla3, 1)
```

```
tipotuberculosis
otra resp
  522 4853
```

Para calcular los marginales por columna

```
> margin.table(tabla3, 2)
```

```
genero
  f    m
1571 3804
```

Para pedir el resumen de la tabla

```
> summary(tabla3)
```

```
Call: xtabs(formula = conteos ~ tipotuberculosis + genero, data = TABLA)
```

```
Number of cases in table: 5375
```

```
Number of factors: 2
```

```
Test for independence of all factors:
```

```
Chisq = 101.41, df = 1, p-value = 7.483e-24
```

Con estos valores se concluye que las variables no son independientes, es decir que la proporción de hombres que muere por tuberculosis de tipo respiratorio  $3534/3804=.929$  es significativamente diferente de la proporción  $1319/1571=.840$  de mujeres que mueren por ese tipo de tuberculosis.

ENCONTRAR ASOCIACION A TRAVES DE LA JI CUADRADA NO IMPLICA NECESARIAMENTE NINGUNA RELACION CAUSAL.

## 2.6 Más de la Ji Cuadrada

### 2.6.1 Sensible al tamaño de muestra

```
> abortoafavor <- c("si", "si", "no", "no")
```

```
> raza <- c("b", "n", "b", "n")
```

```
> conteos <- c(49, 51, 51, 49)
```

```
> TABLA <- data.frame(abortoafavor, raza, conteos)
```

```
> tabla4 <- xtabs(conteos ~ abortoafavor + raza, data = TABLA)
```

```
> summary(tabla4)
```

```
Call: xtabs(formula = conteos ~ abortoafavor + raza, data = TABLA)
```

```
Number of cases in table: 200
```

```
Number of factors: 2
```

Test for independence of all factors:

Chisq = 0.08, df = 1, p-value = 0.7773

```
> abortoafavor <- c("si", "si", "no", "no")
> raza <- c("b", "n", "b", "n")
> conteos <- c(98, 102, 102, 98)
> TABLA <- data.frame(abortoafavor, raza, conteos)
> tabla4 <- xtabs(conteos ~ abortoafavor + raza, data = TABLA)
> summary(tabla4)
```

Call: xtabs(formula = conteos ~ abortoafavor + raza, data = TABLA)

Number of cases in table: 400

Number of factors: 2

Test for independence of all factors:

Chisq = 0.16, df = 1, p-value = 0.6892

```
> abortoafavor <- c("si", "si", "no", "no")
> raza <- c("b", "n", "b", "n")
> conteos <- c(4900, 5100, 5100, 4900)
> TABLA <- data.frame(abortoafavor, raza, conteos)
> tabla4 <- xtabs(conteos ~ abortoafavor + raza, data = TABLA)
> summary(tabla4)
```

Call: xtabs(formula = conteos ~ abortoafavor + raza, data = TABLA)

Number of cases in table: 20000

Number of factors: 2

Test for independence of all factors:

Chisq = 8, df = 1, p-value = 0.004678

El valor de la ji cuadrada queda multiplicado por la constante que multiplique las entradas de la tabla. Los grados de libertad de la ji cuadrada no se modifican, y por tanto con una  $n$  muy grande esta prueba resulta significativa.

## 2.7 Comparación de proporciones en tablas de 2 por 2

Notación para una tabla de 2 por 2

$n_{11}$	$n_{12}$	$n_{1.}$
$n_{21}$	$n_{22}$	$n_{2.}$
$n_{.1}$	$n_{.2}$	$n_{..}$

La estimación de las proporciones está dada por:

$$\hat{\pi}_1 = p_1 = n_{11}/n_{1.} \text{ y } \hat{\pi}_2 = p_2 = n_{21}/n_{2.}$$

## 2.8 Diferencia de proporciones

Para los sujetos en la primera fila se tiene que la probabilidad de éxito es  $\pi_1$  y para la fila 2 es  $\pi_2$ , si comparamos  $\pi_1 - \pi_2$ , las versiones muestrales son  $p_1$  y  $p_2$  cuando las muestras son de tamaño  $N_1$  y  $N_2$ , del curso propedéutico recordamos que para muestras grandes se tiene que:

$$\hat{\sigma}_{p_1-p_2} = \sqrt{\frac{p_1(1-p_1)}{N_1} + \frac{p_2(1-p_2)}{N_2}}$$

Y el intervalo de  $(1 - \alpha)\%$  de confianza para  $\pi_1 - \pi_2$  es

$$(p_1 - p_2) \pm z_{\alpha/2} \hat{\sigma}_{p_1-p_2}$$

## 2.9 Riesgo Relativo

Tal vez ocurre que una diferencia entre proporciones sea más importante cuando se está en los extremos, cerca de 0 o de 1 que cuando se está en el centro. La diferencia entre .010 y .001 es la misma que entre .410 y .401, pero

la primera diferencia es más fuerte pues una es 10 veces la otra, entonces es preferible que consideremos el cociente de proporciones.

En tablas de  $2 \times 2$  el *riesgo relativo* es **el cociente de las probabilidades de éxito en los dos grupos**,

$$\frac{\pi_1}{\pi_2}$$

Para el ejemplo anterior se tiene:

$$r.r_{caso1} = \frac{0.010}{0.001} = 10$$

$$r.r_{caso2} = \frac{.410}{.401} = 1.02$$

OJO hay que definir aquí cual es la variable respuesta para definir la probabilidad de éxito. Es decir qué ponemos como columna y qué como renglón.

Se tiene un riesgo relativo igual a 1 cuando  $\pi_1 = \pi_2$ , es decir cuando la respuesta es independiente del grupo.

## 2.10 Cociente de Momios

Para la fila 1 el *momio* está dado por:  $\text{momio}_1 = \frac{\pi_1}{(1-\pi_1)}$  y para la fila 2 por  $\text{momio}_2 = \frac{\pi_2}{(1-\pi_2)}$ . Así si  $\pi_1 = .75$  entonces el momio es  $.75/.25=3$ . Entonces si el momio=4, el éxito es cuatro veces más probable que un fracaso. Esperamos ver 4 éxitos por cada fracaso.

Si despejamos la probabilidad se tiene que  $\pi = \frac{\text{momio}}{\text{momio}+1}$ , si momio=4 entonces  $\pi = 4/(4+1) = .8$

Ahora el cociente de momios se define como:

$$\theta = \frac{\text{momio}_1}{\text{momio}_2} = \frac{\frac{\pi_1}{(1-\pi_1)}}{\frac{\pi_2}{(1-\pi_2)}}$$

Este NO es un cociente de probabilidades como en el riesgo relativo.

Si  $X$  y  $Y$  son independientes  $\pi_1 = \pi_2$ ,  $\text{momio}_1 = \text{momio}_2$  y también  $\theta = \text{momio}_1/\text{momio}_2 = 1$ . Cuando  $1 < \theta < \infty$  los momios de éxito son mayores

en la fila 1 que en la 2, es decir  $\pi_1 > \pi_2$ , cuando  $0 < \theta < 1$  un éxito es menos probable en la fila 1 que en la 2, es decir  $\pi_1 < \pi_2$ .

Cuando  $\theta$  se aleja del 1, ya sea hacia arriba o hacia abajo, representa mayores niveles de asociación. Una  $\theta = 4$  está más lejos de la independencia que una  $\theta = 2$ , lo mismo una  $\theta = .25$  está más lejos de la independencia que una  $\theta = .50$ .

ATENCIÓN!!!!!!!!!!!!!!

Si se intercambian las filas y se tenía una  $\theta = 4$  ahora se tendrá una  $\theta = .25$ , lo mismo ocurre si se voltean las columnas.

La  $\theta$  no cambia si la tabla se presenta traspuesta, es decir las columnas son los renglones y los renglones son columnas. Como las trata de manera simétrica, no importa cual variable es considerada como respuesta. OJO en riesgo relativo si importa.

El estimador de  $\theta$  está dado por

$$\hat{\theta} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$$

Ejemplo de Infarto

```
> infarto <- c("isi", "isi", "no", "no")
> grupo <- c("placebo", "taspirin", "placebo", "taspirin")
> conteos <- c(189, 104, 10845, 10933)
> TABLA <- data.frame(infarto, grupo, conteos)
> tabla4 <- xtabs(conteos ~ grupo + infarto, data = TABLA)
> tabla4
```

	infarto	
grupo	isi	no
placebo	189	10845
taspirin	104	10933



```
> margin.table(tabla4, 1)
```

```
grupo
```

```
  placebo  taspirin
```

```
  11034    11037
```

```
> margin.table(tabla4, 2)
```

```
infarto
```

```
  isi    no
```

```
  293 21778
```

```
> summary(tabla4)
```

```
Call: xtabs(formula = conteos ~ grupo + infarto, data = TABLA)
```

```
Number of cases in table: 22071
```

```
Number of factors: 2
```

```
Test for independence of all factors:
```

```
  Chisq = 25.014, df = 1, p-value = 5.692e-07
```

Claramente se ve que no son independientes, pero para dónde están jalando las cosas.

Momio para el grupo con placebo:

```
> 189/10845
```

```
[1] 0.01742739
```

Momio pra el grupo con aspirina

```
> 104/10933
```

```
[1] 0.009512485
```

Cálculo de Cociente de momios:

```
> (189/10845)/(104/10933)
```

```
[1] 1.832054
```

Los momios pues son 83% más grandes para el grupo placebo.

## 2.11 Residuales

Una forma de ver qué categorías son las que provocan la asociación de las variables es fijarse en aquellas que tengan los más grandes *residuales ajustados* dados por:

$$\frac{(n_{ij} - E_{ij})}{\sqrt{E_{ij} * (1 - p_{i.})(1 - p_{.j})}}$$

## 2.12 Ejemplo

Ejemplo con una tabla de 2x3.

```
> partido <- c("democ", "democ", "independ", "independ", "republic", "republ.
> genero <- c("fem", "masc", "fem", "masc", "fem", "masc")
> conteos <- c(279, 165, 73, 47, 225, 191)
> TABLA <- data.frame(partido, genero, conteos)
> TABLA <- xtabs(conteos ~ genero + partido, data = TABLA)
> TABLA
```

	partido		
genero	democ	independ	republic
fem	279	73	225
masc	165	47	191

Para estos datos calculamos las marginales por renglon y columna y el total.

```
> margin.table(TABLA, 1)
```

```
genero
```

```
  fem masc
```

```
577 403
```

```
> margin.table(TABLA, 2)
```

```
partido
```

```
  democ independ republic
```

```
  444      120      416
```

```
> sum(TABLA)
```

```
[1] 980
```

Se calcula la prueba ji cuadrada, y podemos acceder a los valores esperados y los residuales como se muestra

```
> prueba <- chisq.test(xtabs(conteos ~ genero + partido, data = TABLA))
```

```
> prueba
```

```
      Pearson s Chi-squared test
```

```
data:  xtabs(conteos ~ genero + partido, data = TABLA)
```

```
X-squared = 7.0095, df = 2, p-value = 0.03005
```

```
> prueba$expected
```

```
      partido
```

```
genero  democ independ republic
```

```
  fem 261.4163 70.65306 244.9306
```

```
  masc 182.5837 49.34694 171.0694
```

```
> prueba$p.value
```

```
[1] 0.03005363
```

```
> prueba$residuals
```

```
      partido
genero  democ  independ  republic
fem    1.0875350  0.2792134 -1.2735005
masc  -1.3013036 -0.3340963  1.5238229
```

Para calcular los residuales ajustados debemos hacerle un pequeño ajuste

```
> proprenqlon <- (margin.table(TABLA, 1)/sum(TABLA))
```

```
> proprenqlon
```

```
genero
      fem      masc
0.5887755 0.4112245
```

```
> propcol <- (margin.table(TABLA, 2)/sum(TABLA))
```

```
> propcol
```

```
partido
      democ  independ  republic
0.4530612 0.1224490 0.4244898
```

```
> auxiliar <- as.matrix(1 - proprenqlon) %*% t(1 - as.matrix(propcol))
```

```
> resajustados <- matrix(nr = 2, ncol = 3)
```

```
> for (i in 1:2) {
```

```
+   for (j in 1:3) {
```

```
+       resajustados[i, j] <- (prueba$residuals[i, j])/sqrt(auxiliar[i, j])
```

```
+   }
```

```
+ }
```

```

> resajustados

           [,1]      [,2]      [,3]
[1,]  2.293160  0.4647941 -2.61778
[2,] -2.293160 -0.4647941  2.61778

> p1x <- 577/980
> thetademrepublic <- (279 * 191)/(225 * 165)
> thetademrepublic

[1] 1.435394

```

Los residuales ajustados se muestran grandes en las mujeres demócratas y en los hombres republicanos, esto se muestra también en la  $\theta = 1.43$ , entonces decimos que los momios de identificarse con los demócratas en vez de los republicanos son 44% más grandes en las mujeres que en los hombres.

## 2.13 Intervalo de confianza para $\log(\theta)$

La distribución muestral del riesgo relativo y cociente de momios es muy asimétrica, debido a esto se usa la función  $\log(\theta)$ . Esta función resulta simétrica alrededor del cero en el sentido de que si se invierte el orden de las columnas o renglones entonces por ejemplo  $\log(2.0) = 0.7$  (al voltear los renglones) se tendría  $\log(0.5) = -0.7$ , digamos que representa el mismo nivel de asociación. El doblar el logaritmo de cocientes de momios representa elevar al cuadrado el cociente de momios. La distribución de  $\log(\theta)$  sigue siendo asimétrica pero muchos más cercana a la normal.

Para una muestra grande tiene media  $\log(\theta)$  y una desviación estándar asintótica de:

$$ASE(\log(\hat{\theta})) = \sqrt{1/n_{11} + 1/n_{12} + 1/n_{21} + 1/n_{22}}$$

Entonces el intervalo de confianza para  $\log(\theta)$  es de la forma:

$$\log(\hat{\theta}) \pm z_{\alpha/2} ASE(\log(\hat{\theta}))$$

Si se exponencian los extremos del intervalo se obtiene un intervalo de confianza para  $\theta$

Para calcular el intervalo de confianza para  $\log(\theta)$  en el ejemplo anterior hacemos:

```
> ase <- sqrt(1/189 + 1/10933 + 1/10845 + 1/104)
```

```
> ase
```

```
[1] 0.1228416
```

```
> log(1.823) - ase * qnorm(0.975)
```

```
[1] 0.3597183
```

```
> log(1.823) + ase * qnorm(0.975)
```

```
[1] 0.8412487
```

```
> exp(log(1.823) - ase * qnorm(0.975))
```

```
[1] 1.432926
```

```
> exp(log(1.823) + ase * qnorm(0.975))
```

```
[1] 2.319261
```

Entonces, el intervalo [1.43,2.31] no contiene al 1, los momios son diferentes para cada grupo, viendo el límite inferior del intervalo se tiene que LOS MOMIOS DE INFARTO AL MIOCARDIO SON AL MENOS 43% MAS ALTOS EN EL GRUPO DE PLACEBO QUE EN EL GRUPO DE ASPIRINA.

El intervalo NO es simétrico. Si una celda es cero la  $\hat{\theta}$  está indefinida o es cero. Si se usa este otro estimador no se tendría ese problema:

$$\tilde{\theta} = \frac{(n_{11} + .5)(n_{22} + .5)}{(n_{12} + .5)(n_{21} + .5)}$$

con desviación estándar asintótica de

$$ASE(\tilde{\theta}) = \sqrt{1/(n_{11} + .5) + 1/(n_{12} + .5) + 1/(n_{21} + .5) + 1/(n_{22} + .5)}$$

En el ejemplo anterior  $\tilde{\theta} = 1.828$  muy cercano a  $\hat{\theta} = 1.832$

## 2.14 Intervalo de confianza para $\log(RR)$

También el riesgo relativo tiene una distribución muy asimétrica y de manera análoga el intervalo de confianza para  $\log(RR)$  está dado por:

$$\log(\hat{RR}) \pm z_{\alpha/2} * \sqrt{\frac{1 - p_1}{n1. * p_1} + \frac{1 - p_2}{n2. * p_2}}$$

### 2.14.1 Ejemplo caso y controles

Los datos se refieren a 262 mujeres de edad intermedia (menores a 69 años) que son admitidas en unidades médicas con infarto agudo al miocardio (MI) en un lapso de 5 años, cada caso es apareado con dos pacientes control recibidos en los mismos hospitales con algún otro padecimiento agudo. Se les clasifica como *si*, sin son fumadoras o exfumadoras y como *no* a aquellas que nunca han fumado. Por el diseño la distribución marginal de MI, está fija, habiendo 2 controles por cada caso. Estos estudios son conocidos como *caso y controles*, este diseño permite tener suficientes casos con la enfermedad (característica) de interés, y después se buscan ciertas características en su historia clínica, es decir es un estudio retrospectivo.

Se desea comparar fumadoras versus no fumadoras en cuanto a la proporción de personas que sufren infarto al miocardio. Esto se refiere a la

distribución marginal de MI dado el estatus de fumador. En esta muestra aproximadamente un tercio de ésta sufrió MI, no tiene sentido usar  $1/3$  como estimador de la probabilidad de sufrir MI. ( $1/3 = P(MI) = P(MI/F)P(F) + P(MI/NF)P(NF)$ ). Lo que si se puede calcular es la distribución condicional de ser fumador dado que se sufrió un MI.

```

      infarto
fuma control miocardio
NO      346      90
SI      173     172

```

Pearsons Chi-squared test with Yates continuity correction

```

data: xtabs(conteos ~ fuma + infarto, data = TABLA)
X-squared = 72.4241, df = 1, p-value < 2.2e-16

```

El cociente momios es:

```

> theta <- (172 * 346)/(90 * 173)
> theta

[1] 3.822222

```

Para las mujeres que sufrieron infarto, la proporción de fumadoras es:

```

> p1 <- 172/262
> p1

[1] 0.6564885

```

Para las mujeres que no sufrieron infarto, la proporción de fumadoras es:

```

> p2 <- 173/519
> p2

```



[1] 0.3333333

Como vimos que la probabilidad de sufrir infarto es pequeña para ambos grupos, entonces podemos pensar que el  $r.r$  es parecido a 3.82. Entonces decimos que: *las mujeres que han fumado alguna vez tienen una probabilidad de sufrir un infarto casi 4 (3.82) veces mayor que las mujeres que no han fumado.*