

# Datos Categóricos

## Clase 9

Juan Carlos Correa

4 de abril de 2022

**Resultado Importante I** *Suponga que  $X_n$  es  $AN(\mu, \sigma_n^2)$  con  $\sigma_n \rightarrow 0$ . Sea  $g$  una función de valor real diferenciable en  $X = \mu$  con  $g'(\mu) \neq 0$ . Entonces*

$$g(X_n) \sim AN\left(g(\mu), [g'(\mu)]^2 \sigma_n^2\right)$$

**Resultado Importante II** Sea  $\mathbf{X}_n = (X_{n1}, X_{n2}, \dots, X_{nk})'$  y además asuma que

$$\mathbf{X}_n \sim AN(\mu, b_n^2 \Sigma)$$

con  $\Sigma$  matriz de covarianzas y  $b_n \rightarrow 0$ .

Sea  $g(\mathbf{x}) = (g_1(\mathbf{x}), g_2(\mathbf{x}), \dots, g_m(\mathbf{x}), )'$ , donde  $\mathbf{x} = (x_1, x_2, \dots, x_k)'$ , una función con argumento un vector y donde cada componente es una función de valor real y tiene un diferencial no cero  $g_i(\mu; \mathbf{t})$ ,  $\mathbf{x} = (t_1, t_2, \dots, t_k)'$ , en  $\mathbf{x} = \mu$ . Haga

$$\mathbf{D} = \left[ \frac{\partial g_i}{\partial x_j} \Big|_{\mathbf{x}=\mu} \right]_{m \times k}$$

Entonces  $g(\mathbf{X}_n) \sim AN(g(\mu), b_n^2 \mathbf{D} \Sigma \mathbf{D}')$

## La Razón de Odds

El estimador muestral de  $\psi$  será

$$r = \frac{\left( \frac{n_{11}}{n_{+1}} \right)}{\left( \frac{n_{12}}{n_{+2}} \right)} = \frac{\frac{n_{11}}{n_{21}}}{\frac{n_{12}}{n_{22}}} = \frac{n_{11}n_{22}}{n_{12}n_{21}}$$

para lo anterior, se presupone una tabla conteos de como la que aparece a continuación

	$A$	$A^c$	
$B$	$n_{11}$	$n_{12}$	$n_{1+}$
$B^c$	$n_{21}$	$n_{22}$	$n_{2+}$
	$n_{+1}$	$n_{+2}$	

**Problema con celdas con ceros** Un problema con este estimador  $r$  es la presencia de ceros en las celdas, ya que puede convertirse en una forma indeterminada.

Varios estimadores adicionales han sido propuestos para la razón odds y para el logarítmico de la razón de odds. Entre ellos tenemos:

- El de Haldane:

$$\hat{\psi}_H = \frac{(a + \frac{1}{2})(d + \frac{1}{2})}{(c + \frac{1}{2})(b + \frac{1}{2})}$$

- El de Jewell:

$$\hat{\psi}_J = \frac{ad}{(b + 1)(c + 1)}$$

- Estimador de máxima verosimilitud condicional: Este estimador es la solución a un polinomio de alto grado de la forma:

$$\sum_{j=s}^{\delta} \binom{N_1}{j} \binom{N_2}{k_1 - j} (a - j) \rho^j$$

donde  $s = \max(0, k_1 - N_2)$  y  $\delta = \min(k_1, N_1)$

**Propiedades de la razón de odds** Algunas propiedades de la razón de odds son las siguientes:

- Es un número nonegativo.
- Cuando todas las celdas tienen probabilidades positivas, la independencia entre las dos variables es equivalente a  $\psi = 0$ .
- Es invariante bajo el intercambio de filas o columnas.
- Es invariante bajo multiplicaciones de filas y columnas.

- La interpretación es clara. Valores de  $\psi$  que se alejen de 1.0 en una dirección particular representa una asociación fuerte. Dos valores de  $\psi$  pueden representar un mismo nivel de asociación (un valor y su inverso) pero en direcciones opuestas. Para simetrizar esta medida se trabaja con el  $\log(\psi)$ . Valores menores que uno indican una asociación negativa, mientras valores mayores que 1 indican una asociación positiva.
- Puede usarse en tablas  $I \times J$  (y tablas multidimensionales) mirando series de particiones  $2 \times 2$  o mirando subtablas  $2 \times 2$ .



## Distribución asintótica de la Razón de Odds: Esquema de muestreo *multinomial*

Sean

$$(n_1, \dots, n_k) \sim \text{Multinomial}(\pi, n)$$

$$\pi = (\pi_1, \pi_2, \dots, \pi_k)^T$$

$$n = n_1 + \dots + n_k$$

Una estimación para el vector  $\pi$  es el vector

$$\hat{\pi} = (\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_k)^T.$$

La  $i$ -ésima observación es

$$\mathbf{Y}_i = (Y_{i1}, Y_{i2}, \dots, Y_{ik})'$$

donde

$$Y_{ij} = \begin{cases} 1 & \text{si cae en la celda } j \\ 0 & \text{en otro caso} \end{cases}$$

y además

$$\sum_j Y_{ij} = 1$$

Ahora

$$E[\mathbf{Y}_i] = \boldsymbol{\pi}$$

$$\text{cov}(\mathbf{Y}_i) = \boldsymbol{\Sigma} \quad i = 1, \dots, n$$

$$\sigma_{jj} = \text{var}(Y_{ij}) = \pi_j(1 - \pi_j)$$

$$\sigma_{jk} = \text{cov}(Y_{ij}, Y_{ik}) = E(Y_{ij}Y_{ik}) - E(Y_{ij})E(Y_{ik})$$

$$= -\pi_j\pi_k \quad j \neq k$$

$$\boldsymbol{\Sigma} = \text{Diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}^T$$

$$\hat{\boldsymbol{\pi}} = \frac{1}{n} \sum_{i=1}^n \mathbf{Y}_i$$

$$\text{cov}(\hat{\boldsymbol{\pi}}) = \frac{(\text{Diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}^T)}{n} \rightarrow \text{Matriz singular}$$

**Teorema central del límite multivariable** Bajo el supuesto que  $\mathbf{Y}_i, i = 1, \dots, n$  sea una muestra aleatoria de una distribución  $Multinomial(\pi, 1)$ , entonces

$$\sqrt{n}(\hat{\pi} - \pi) \xrightarrow{a} N(\mathbf{0}, \text{Diag}(\pi) - \pi\pi^T)$$

cuando  $n \rightarrow \infty$ .

Ahora

$$\begin{aligned}g(\pi) &= \log(\pi) \\ \frac{\partial g}{\partial \pi} &= \text{Diag}(\pi)^{-1}\end{aligned}$$

La covarianza de la matriz asintótica de

$$\sqrt{n} [\log(\hat{\pi}) - \log(\pi)]$$

es

$$\text{Diag}(\pi)^{-1} [\text{Diag}(\pi) - \pi\pi^T] \text{Diag}(\pi)^{-1} = \text{Diag}(\pi)^{-1} - \mathbf{1}\mathbf{1}^T$$

Para una matriz  $C$  de constantes

$$\sqrt{n}C [\log(\hat{\pi}) - \log(\pi)] \xrightarrow{a} N\left(\mathbf{0}, C\text{Diag}(\pi)^{-1}C^T - C\mathbf{1}\mathbf{1}^T C^T\right)$$

Con base en el anterior resultado, consideremos el siguiente vector

$$\begin{pmatrix} \pi_{11} \\ \pi_{12} \\ \pi_{21} \\ \pi_{22} \end{pmatrix}$$

El Odds ratio será

$$OR = \psi = \frac{\frac{\pi_{11}}{\pi_{21}}}{\frac{\pi_{12}}{\pi_{22}}} = \frac{\pi_{11}\pi_{22}}{\pi_{12}\pi_{21}}$$

Ahora

$$\log(\psi) = C(\log(\pi)) = [1 \quad -1 \quad -1 \quad 1] \begin{bmatrix} \pi_{11} \\ \pi_{12} \\ \pi_{21} \\ \pi_{22} \end{bmatrix}$$

entonces

$$\begin{aligned} nVar(\log(\hat{\psi})) &= C \text{Diag}(\pi)^{-1} C^T - C \mathbf{1} \mathbf{1}^T C^T \\ &= [1 \quad -1 \quad -1 \quad 1] \begin{bmatrix} \frac{1}{\pi_{11}} & 0 & 0 & 0 \\ 0 & \frac{1}{\pi_{12}} & 0 & 0 \\ 0 & 0 & \frac{1}{\pi_{21}} & 0 \\ 0 & 0 & 0 & \frac{1}{\pi_{22}} \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} \\ &\quad - [1 \quad -1 \quad -1 \quad 1] \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} [1 \quad 1 \quad 1 \quad 1] \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} \end{aligned}$$

$$\begin{aligned}
&= \begin{bmatrix} \frac{1}{\pi_{11}} & -\frac{1}{\pi_{12}} & -\frac{1}{\pi_{21}} & \frac{1}{\pi_{22}} \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} \\
&= \frac{1}{\pi_{11}} + \frac{1}{\pi_{12}} + \frac{1}{\pi_{21}} + \frac{1}{\pi_{22}}
\end{aligned}$$



### Distribución Asintótica de $\log(\hat{\psi})$

$$\log(\hat{\psi}) \sim AN\left(\log(\psi), \frac{1}{n} \left( \frac{1}{\pi_{11}} + \frac{1}{\pi_{12}} + \frac{1}{\pi_{21}} + \frac{1}{\pi_{22}} \right)\right)$$

Un intervalo de confianza para  $\log(\psi)$  del 95 % es

$$\left( \log(\hat{\psi}) \mp 1,96 \sqrt{\frac{1}{n} \left( \frac{1}{\hat{\pi}_{11}} + \frac{1}{\hat{\pi}_{12}} + \frac{1}{\hat{\pi}_{21}} + \frac{1}{\hat{\pi}_{22}} \right)} \right)$$

o

$$\left( \log(\hat{\psi}) \mp 1,96 \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}} \right)$$

Una forma muy común de hallar un intervalo de confianza para  $\psi$  se calcula invirtiendo el intervalo anterior

$$LI = \exp \left( \log (\hat{\psi}) - 1,96 \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}} \right)$$

y

$$LS = \exp \left( \log (\hat{\psi}) + 1,96 \sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}} \right)$$

## Distribución Asintótica de $\hat{\psi}$

Tenemos

$$Y = \log(\hat{\psi}) \sim AN\left(\log(\psi), \frac{1}{n} \left( \frac{1}{\pi_{11}} + \frac{1}{\pi_{12}} + \frac{1}{\pi_{21}} + \frac{1}{\pi_{22}} \right)\right)$$

y que

$$\hat{\psi} = \exp(\log(\hat{\psi})) = e^Y$$

Entonces

$$Y \sim AN\left(\exp(\log(\psi)), \frac{1}{n} (\exp(\log(\psi)))^2 \left( \frac{1}{\pi_{11}} + \frac{1}{\pi_{12}} + \frac{1}{\pi_{21}} + \frac{1}{\pi_{22}} \right)\right)$$

$$Y \sim AN\left(\psi, \psi^2 \left( \frac{1}{\pi_{11}} + \frac{1}{\pi_{12}} + \frac{1}{\pi_{21}} + \frac{1}{\pi_{22}} \right)\right)$$

Un intervalo de confianza para  $\psi$  del 95 % basado en la distribución anterior es

$$LI = \hat{\psi} - 1,96\hat{\psi}\sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}$$

y

$$LS = \hat{\psi} + 1,96\hat{\psi}\sqrt{\frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}}$$

## Programa en *R* para calcular la razón de odds

### Intervalo aproximado

```
> intervalo.razon.odds<-function(Tabla,nivel=0.95,correccion=0.5){  
  Tabla<-ifelse(Tabla==0,0.5,Tabla)  
  odds<-Tabla[1,1]*Tabla[2,2]/(Tabla[1,2]*Tabla[2,1])  
  error<-odds*sqrt(1/Tabla[1,1]+1/Tabla[1,2]+1/Tabla[2,1]+1/Tabla[2,2])  
  z<-qnorm(0.5+nivel/2)  
  LI<-odds-z*error  
  LS<-odds+z*error  
  list(odds=odds,error=error,LI=LI,LS=LS)  
}
```

```
>nacimientos.medellin<-matrix(c(4757,430,5148,464),ncol=2,byrow=T)
> nacimientos.medellin
      [,1] [,2]
[1,] 4757  430
[2,] 5148  464
> intervalo.razon.odds(nacimientos.medellin)
$odds
[1] 0.9971124

$error
[1] 0.06969253

$LI
[1] 0.8605176

$LS
[1] 1.133707

>
> odds.nacimientos<-intervalo.razon.odds(nacimientos.medellin)
> odds.nacimientos$LI
```

```
[1] 0.8605176
> odds.nacimientos$LS
[1] 1.133707
> odds.nacimientos$odds
[1] 0.9971124
>
```

```
> fisher.test(matrix(c(4757,430,5148,464),ncol=2,byrow=T))
```

Fisher's Exact Test for Count Data

```
data: matrix(c(4757, 430, 5148, 464), ncol = 2, byrow = T)
```

```
p-value = 0.9721
```

```
alternative hypothesis: true odds ratio is not equal to 1
```

```
95 percent confidence interval:
```

```
0.8674057 1.1463830
```

```
sample estimates:
```

```
odds ratio
```

```
0.9971126
```

```
>
```



## Cálculos en R

Podemos usar la librería *Epi*

```
> library(Epi)
> twoby2(matrix(c(4757, 430, 5148, 464),ncol=2,byrow=T))
2 by 2 table analysis:
```

```
-----
Outcome      : Col 1
Comparing    : Row 1 vs. Row 2
```

	Col 1	Col 2	P(Col 1)	95% conf. interval	
Row 1	4757	430	0.9171	0.9093	0.9243
Row 2	5148	464	0.9173	0.9098	0.9242

	95% conf. interval		
Relative Risk:	0.9998	0.9885	1.0112
Sample Odds Ratio:	0.9971	0.8695	1.1435
Probability difference:	-0.0002	-0.0122	0.0117

Asymptotic P-value: 0.967

```
-----
```