

Modelo Lineal Generalizado

Notemos que el estadístico χ^2 puede escribirse como

$$\chi^2 = \mathbf{e}'\mathbf{e}$$

donde

$$\mathbf{e}' = \left(\sqrt{n} \frac{p_1 - \pi_1(\bar{\theta})}{\sqrt{\pi_1(\bar{\theta})}}, \dots, \sqrt{n} \frac{p_N - \pi_N(\bar{\theta})}{\sqrt{\pi_N(\bar{\theta})}} \right)$$

Para derivar la distribución asintótica de χ^2 necesitaremos la conjunta de $(\mathbf{p}, \boldsymbol{\pi}(\bar{\theta}))$ y deduciremos que

$$\mathbf{e} \xrightarrow{\mathcal{D}} N(0, I - \boldsymbol{\pi}(\theta_0)^{1/2} \boldsymbol{\pi}'(\theta_0)^{1/2} - A(A'A)^{-1}A')$$

Teorema: Sea Y un vector con distribución $N(\boldsymbol{\nu}, \Sigma)$. Una condición necesaria y suficiente para que $(Y - \boldsymbol{\nu})' \mathbf{C} (Y - \boldsymbol{\nu})$ tenga distribución χ^2 es que $\Sigma \mathbf{C} \Sigma \mathbf{C} \Sigma = \Sigma \mathbf{C} \Sigma$, donde los grados de libertad serán el rango de $\mathbf{C} \Sigma$ (si Σ es no singular la condición se simplifica a $\mathbf{C} \Sigma \mathbf{C} = \mathbf{C}$). (Rao, 1965, p. 150)

Como hemos visto, $\chi^2 = \mathbf{e}' \mathbf{e}$, luego aplicaremos el resultado de Rao con $\boldsymbol{\nu} = 0$, $\mathbf{C} = I$, $\Sigma = I - \boldsymbol{\pi}(\theta_0)^{1/2} \boldsymbol{\pi}'(\theta_0)^{1/2} - A(A'A)^{-1}A'$, por lo que resultará que

$$\mathbf{e}' \mathbf{e} \xrightarrow{\mathcal{D}} \chi_{N-1-q}^2$$

Volviendo al Test de Independencia:

En una tabla de $I \times J$ con muestreo multinomial, la hipótesis nula de independencia equivale a

$$H_0 : \pi_{ij} = \pi_{i+}\pi_{+j} \quad \forall i, j$$

Usando el estadístico de Pearson tendríamos

$$\sum_{i=1}^I \sum_{j=1}^J \frac{(n_{ij} - \widehat{m}_{ij})^2}{\widehat{m}_{ij}}$$

Respecto a los grados de libertad, estos están determinados por la cantidad de casillas y de parámetros, que en este caso serán

$$I * J - 1 - (I - 1) - (J - 1) = (I - 1) * (J - 1).$$

Volvamos al ejemplo de la **filiación partidaria** que vimos en la primera clase. En la siguiente tabla tenemos los valores observados y en rojo los valores predichos bajo el modelo de independencia

S: Sexo	C: Identificación partidaria			Total
	Demócrata	Independiente	Republicano	
Mujer	279 (261.4)	73 (70.7)	225 (244.9)	577
Hombre	165 (182.6)	47 (49.3)	191 (171.1)	403
Total	444	120	416	980

Cuadro 12: Datos de la General Social Survey, 1991.

El valor del estadístico test de χ^2 es 7.01, con $IJ - 1 - (I - 1) - (J - 1) = (I - 1)(J - 1) = (2 - 1)(3 - 1) = 2$ grados de libertad. El p-valor correspondiente es 0.03, de manera que a los valores habituales se rechazaría la hipótesis de independencia, indicando que el sexo y la identificación partidaria estarían asociados.

Otro Ejemplo

Este es otro ejemplo en que las probabilidades dependen de una cantidad menor de parámetros desconocidos, θ .

Una muestra de 156 terneros nacidos en Florida fueron clasificados de acuerdo a que hayan contraído neumonía dentro de los 60 días de haber nacido. Los terneros que contrayeron neumonía fueron a su vez clasificados en si se reinfectaron o no a los 15 días de haberse curado. La Tabla muestra los datos recolectados:

	Segunda Infección	
	Si	No
Primera Infección		
Si	30	63
No	0	63

Cuadro 13:

Es claro que los terneros que no tuvieron una primera infección no pudieron reinfectarse, es por ello que ninguna observación puede verse en la casilla 21 y en consecuencia en la tabla $n_{21} = 0$. Esto es lo que se conoce como un **cero**

estructural. El objetivo en este estudio era testear si la probabilidad de una primera infección era igual que la probabilidad de una segunda infección, dado que el ternero había contraído una primera infección.

Es decir, la hipótesis a testear es

$$H_0 : \pi_{11} + \pi_{12} = \frac{\pi_{11}}{\pi_{11} + \pi_{12}}$$

o equivalentemente $\pi_{11} = (\pi_{11} + \pi_{12})^2$. De manera que si llamamos $\pi = \pi_{11} + \pi_{12}$ a la probabilidad de infección primaria, el modelo bajo H_0 corresponde a una **trinomial** como muestra la siguiente tabla:

	Segunda Infección		
	Si	No	Total
Primera Infección			
Si	π^2	$\pi(1 - \pi)$	π
No	–	$1 - \pi$	$1 - \pi$

Cuadro 14:

En este caso el likelihood resulta

$$(\pi^2)^{n_{11}}(\pi(1 - \pi))^{n_{12}}(1 - \pi)^{n_{22}},$$

el log-likelihood queda

$$n_{11} \log(\pi^2) + n_{12} \log(\pi(1 - \pi)) + n_{22} \log(1 - \pi),$$

Derivando e igualando a 0 resulta

$$\hat{\pi} = \frac{2n_{11} + n_{12}}{2n_{11} + 2n_{12} + n_{22}}$$

En la siguiente tabla se muestran en rojo (2do. renglón) los valores esperados bajo H_0

El estadístico de Pearson da $\chi^2 = 19.7$ con un total de $(3-1)-1=1$ grados de libertad. Dado que el p-valor es 0.00001 hay una fuerte evidencia contra H_0 . Si miramos la tabla encontramos que muchos más terneros contraen una primera infección y no la segunda de lo que el modelo bajo H_0 predice. Con esto los investigadores concluyeron que la primera infección tiene un efecto inmunizador.

	Segunda Infección	
	Si	No
Primera Infección		
Si	30 (38.1)	63 (39)
No	0 (-)	63 (78.9)

Cuadro 15:

Estadístico G^2

Otra medida alternativa para la distancia entre $\hat{\boldsymbol{\pi}}$ y \mathbf{p} muy usada es la **deviance** G^2 , que es un estadístico basado en el cociente de verosimilitud.

Si queremos testear

H_0 : Modelo restringido ω

H_1 : Modelo Saturado Ω ,

el cociente estaría dado por

$$\Lambda = \frac{\text{máx}_\omega L}{\text{máx}_\Omega L}$$

Si consideramos $G^2 = -2 \log \Lambda$ queda definido el estadístico como

$$\begin{aligned} G^2 &= -2 \log \Lambda = 2[l(\mathbf{p}, \mathbf{X}) - l(\widehat{\boldsymbol{\pi}}, \mathbf{X})] \\ &= 2 \left[\sum_{i=1}^N X_i \log p_i - \sum_{i=1}^N X_i \log \widehat{\pi}_i \right] \\ &= 2 \left[\sum_{i=1}^N X_i \log \frac{p_i}{\widehat{\pi}_i} \right] \\ &= 2 \left[\sum_{i=1}^N X_i \log \frac{X_i}{n \widehat{\pi}_i} \right] \end{aligned}$$

Probaremos que bajo H_0 la distribución límite de G^2 es también χ^2 con $N - 1 - \#$ parámetros bajo ω , es decir la misma distribución límite que la del estadístico de Pearson.

Para derivar la distribución asintótica, probaremos que $G^2 - \chi^2 \xrightarrow{p} 0$.

Una ventaja de G^2 es que tiene sentido en modelos más generales.

En el ejemplo de **Identificación Partidaria vs. Sexo**, $G^2 = 7$, que da también un p-valor de 0.03.

Efecto de observar ceros

Si en alguna celda se observa un 0, el estadístico χ^2 puede calcularse sin problemas, siempre que las $\hat{\pi}$'s sean todas positivas. Sin embargo, el estadístico G^2 tiene problemas, pues si $X_i = 0$, entonces $X_i \log \frac{X_i}{n\hat{\pi}_i}$ no está definido. Si reescribimos a G^2 como

$$\begin{aligned} G^2 &= -2 \log \Lambda = 2 \log \frac{L(\mathbf{p}, \mathbf{X})}{L(\hat{\boldsymbol{\pi}}, \mathbf{X})} \\ &= 2 \log \prod_{i=1}^N \left(\frac{X_i/n}{\hat{\pi}_i} \right)^{X_i} \end{aligned}$$

es claro que una celda con un 0 aporta un 1 al producto y por lo tanto puede ser ignorada. Luego podemos calcular a G^2 como

$$2 \sum_{i: x_i > 0} X_i \log \frac{X_i}{n\hat{\pi}_i}$$

Si alguna $\hat{\pi}$ es 0, los dos estadísticos se *rompen*.

Cuán grande debe ser n para tener una buena aproximación?

Sabemos que a medida que n crece la distribución de χ^2 y de G^2 se aproximan a una distribución límite χ^2 , sin embargo nos preguntamos cuán grande es grande.

- Una vieja regla conocida para las binomiales dice que la aproximación χ^2 es buena si $n\hat{\pi}_i \geq 5$, $i = 1, \dots, N$.

- Otra regla más permisiva establece que la aproximación χ^2 es buena si a lo sumo el 20% de las casillas tienen $n\hat{\pi}_i < 5$, $i = 1, \dots, N$ y ninguna casilla tiene $n\hat{\pi}_i < 1$.

- En tablas *sparse* (esparcidas????)(por ejemplo, $n/N < 5$) la aproximación χ^2 es pobre. En realidad, si los datos están distribuidos en la tabla de forma muy desigual, en el sentido de que hay zonas de la tabla que son *sparse*, la aproximación χ^2 también puede ser pobre, aún cuando el n total sea grande.

Hemos probado que los dos estadísticos se aproximan a 0, si el modelo es cierto. Si el modelo no es cierto, ambos crecen, pero la diferencia entre ambos también puede crecer. De manera, que si el modelo tiene un ajuste pobre los dos estadísticos pueden ser grandes y estar lejos uno de otro, i.e., $|\chi^2 - G^2|$ no necesariamente tiende a 0 con n . Aún en esa situación, los correspondientes

p-valores pueden estar cerca de 0 y podemos llegar a la misma conclusión a partir de ellos.

Para ser más precisos, consideremos una sucesión de situaciones $\boldsymbol{\pi}_n$ para las cuales la falta de ajuste disminuye con n , es decir trabajaremos con alternativas contiguas. Supongamos que el modelo bajo la hipótesis nula es $\boldsymbol{\pi}$, pero en realidad

$$\boldsymbol{\pi}_n = \mathbf{f}(\boldsymbol{\theta}) + \boldsymbol{\delta}/\sqrt{n},$$

entonces si $\boldsymbol{\delta} = 0$, el modelo es cierto.

Para estas alternativas contiguas, Mitra (1958) demostró que el estadístico de Pearson tiene distribución asintótica χ^2 no central, con $N - 1 - q$ grados de libertad, con parámetro de no centralidad dado por

$$\lambda = n \sum_{i=1}^N \frac{(\pi_{ni} - f_i(\boldsymbol{\theta}))^2}{f_i(\boldsymbol{\theta})}$$

Observemos que λ tiene la forma del estadístico χ^2 en el que se reemplazó a \mathbf{p} por $\boldsymbol{\pi}_n$ y a $\widehat{\boldsymbol{\pi}}$ por $\mathbf{f}(\boldsymbol{\theta})$. Análogamente, utilizando los mismos reemplazos obtenemos el parámetro de no centralidad de G^2 . Haberman (1974) demostró que

bajo ciertas condiciones χ^2 y G^2 tienen el mismo parámetro de no centralidad, pero éste no es siempre el caso.

Residuos de Pearson y deviance

Como ya hemos visto podemos escribir al estadístico de Pearson como

$$\chi^2 = \sum_{i=1}^N e_i^2$$

donde

$$e_i = \sqrt{n} \frac{p_i - \hat{\pi}_i}{\sqrt{\hat{\pi}_i}}$$

se conoce como el i -ésimo residuo de Pearson.

Estos residuos se comportan de alguna manera como los residuos estandarizados que conocimos en regresión lineal. Es común que se compare a $|e_i|$ con 2, indicándose falta de ajuste en la i -ésima casilla si $|e_i| > 2$. El análisis de estos residuos puede sugerirnos en que sentido los datos se apartan del modelo ajustado.

De la misma forma, la deviance puede interpretarse como una suma de cuadra-

dos de residuos

$$G^2 = \sum_{i=1}^N e_i^2$$

donde

$$e_i = \sqrt{\left| 2X_i \log \frac{X_i}{n\hat{\pi}_i} \right|} \times \text{sgn}(X_i - n\hat{\pi}_i)$$

que se conocen como componentes de la deviance.

Medidas de Asociación

A fin de describir el grado de asociación entre las variables de una tabla de contingencia es frecuente que se usen distintas medidas.

Comenzaremos con tablas de 2×2 , como las que siguen

X	Y		Total	X	Y		Total
	1	2			1	2	
1	π_{11}	π_{12}	π_{1+}	1	n_{11}	n_{12}	n_{1+}
2	π_{21}	π_{22}	π_{2+}	2	n_{21}	n_{22}	n_{2+}
Total	π_{+1}	π_{+2}	1	Total	n_{+1}	n_{+2}	1

Consideremos la siguiente tabla que corresponde a un informe sobre la relación entre el uso de aspirina y el infarto de miocardio realizado por el Physicians Health Study Research Group de Harvard Medical School:

	Infarto de Miocardio		Total
	si	no	
Aspirina	104	10933	11037
Placebo	189	10845	11034

Diferencia de Proporciones o Riesgo Atribuible

Miremos a Y como variable de respuesta y a X como variable explicativa, tal como sería natural en un muestreo de producto multinomial en que

$$n_{11} \sim Bi(n_{1+}, \frac{\pi_{11}}{\pi_{1+}}) \text{ y } n_{21} \sim Bi(n_{2+}, \frac{\pi_{21}}{\pi_{2+}})$$

independientes.

La diferencia de proporciones se define como

$$\begin{aligned} \delta &= P(Y = 1|X = 1) - P(Y = 1|X = 2) \\ &= \frac{\pi_{11}}{\pi_{1+}} - \frac{\pi_{21}}{\pi_{2+}} \\ &= \pi_{1|1} - \pi_{1|2} \end{aligned}$$

Dado que δ es función de los parámetros de $P(Y|X)$, los estimadores de máxima verosimilitud serán los mismos bajo los tres tipos de muestreo que hemos visto.

Podemos estimar a δ como

$$d = \frac{n_{11}}{n_{1+}} - \frac{n_{21}}{n_{2+}}$$

$$= p_{1|1} - p_{1|2}$$

En el ejemplo de Infarto de Miocardio tenemos

$$d = 104/11037 - 189/11034 = 0.0094 - 0.0171 = -0.0077$$

Observemos que

$$E(d) = E(p_{1|1} - p_{1|2}) = \pi_{1|1} - \pi_{1|2}$$

$$V(d) = V(p_{1|1} - p_{1|2}) = \frac{\pi_{1|1}(1 - \pi_{1|1})}{n_{1+}} + \frac{\pi_{1|2}(1 - \pi_{1|2})}{n_{2+}}$$

siendo la última igualdad cierta por la independencia entre las filas.

Si n_{1+} y n_{2+} son grandes, d es aproximadamente normal, es decir

$$\frac{(p_{1|1} - p_{1|2}) - (\pi_{1|1} - \pi_{1|2})}{\sqrt{\frac{\pi_{1|1}(1 - \pi_{1|1})}{n_{1+}} + \frac{\pi_{1|2}(1 - \pi_{1|2})}{n_{2+}}}}$$

es aproximadamente $N(0, 1)$. Por lo tanto haciendo un plug-in para estimar la varianza podemos obtener un intervalo asintótico para δ de nivel $1 - \alpha$ como

$$d \pm z_{\alpha/2} \sqrt{\frac{p_{1|1}(1-p_{1|1})}{n_{1+}} + \frac{p_{1|2}(1-p_{1|2})}{n_{2+}}}$$

$$p_{1|1} - p_{1|2} \pm z_{\alpha/2} \sqrt{\frac{p_{1|1}(1-p_{1|1})}{n_{1+}} + \frac{p_{1|2}(1-p_{1|2})}{n_{2+}}}$$

Riesgo Relativo

Notemos que que la diferencia entre 41% y 40.1% es la misma que entre 1% y 0.1%. Sin embargo, 1% es diez veces 0.1%. Este es un problema de la diferencia de proporciones. Si estamos trabajando con eventos poco frecuentes $\pi_{1|1}$ y $\pi_{1|2}$ serán muy pequeñas y δ será cercano a 0, aún cuando el efecto sea importante, como en el ejemplo. Esto es frecuente en epidemiología en donde la prevalencia de ciertas enfermedades es muy baja.

Esto sugiere la conveniencia de considerar una medida relativa como el **riesgo relativo**

$$RR = \frac{P(Y = 1|X = 1)}{P(Y = 1|X = 2)} = \frac{\pi_{11}/\pi_{1+}}{\pi_{21}/\pi_{2+}}$$

El riesgo relativo es una medida no negativa y un riesgo relativo igual a 1 corresponde a independencia.

RR esta medida es sólo función de $P(Y|X)$, por lo tanto la inferencia que hagamos sobre RR será la misma para los tres muestreos que hemos considerado. La comparación en la otra respuesta da otro riesgo relativo.

El EMV de RR es

$$rr = \frac{n_{11}/n_{1+}}{n_{21}/n_{2+}}$$

En el ejemplo quedaría:

$$rr = \frac{0,0094}{0,0171} = 0,55,$$

esto significa que el riesgo de infarto de miocardio en el grupo tratado con aspirina es aproximadamente la mitad que en grupo que recibió placebo.

Dado que una aproximación normal a su logaritmo es buena suele usarse como medida $\log(RR)$, que se estima por $\log(rr) = \log p_{1|1} - \log p_{1|2}$.

Sabemos que

$$\sqrt{n_{i+}}(p_{1|i} - \pi_{1|i}) \xrightarrow{D} N(0, \pi_{1|i}(1 - \pi_{1|i})),$$

luego usando el método Δ obtenemos que

$$\sqrt{n_{i+}}(\log p_{1|i} - \log \pi_{1|i}) \xrightarrow{D} N\left(0, \frac{(1 - \pi_{1|i})}{\pi_{1|i}}\right).$$

Por la independencia entre las filas, obtenemos que la varianza asintótica de $\log(rr)$ es

$$V(\log(rr)) \simeq \frac{(1 - \pi_{1|1})}{n_{1+}\pi_{1|1}} + \frac{(1 - \pi_{1|2})}{n_{2+}\pi_{1|2}}$$

y se puede estimar por

$$\begin{aligned} \widehat{V}(\log(rr)) &\simeq \frac{(1 - p_{1|1})}{n_{1+}p_{1|1}} + \frac{(1 - p_{1|2})}{n_{2+}p_{1|2}} \\ &\simeq \frac{1}{n_{11}} - \frac{1}{n_{1+}} + \frac{1}{n_{21}} - \frac{1}{n_{2+}} \end{aligned}$$

Un intervalo de nivel asintótico $1 - \alpha$ para $\log(RR)$ es

$$\log(rr) \pm z_{\alpha/2} \sqrt{\widehat{V}(\log(rr))}$$

Como $\log(rr)$ no existe si algún $\pi_{1|i} = 0$ suele usarse

$$\log(\tilde{r}\tilde{r}) = \log\left(\frac{n_{11} + 1/2}{n_{1+} + 1/2}\right) - \log\left(\frac{n_{21} + 1/2}{n_{2+} + 1/2}\right)$$

Odds Ratio (Producto Cruzado)

El riesgo relativo es el cociente de dos probabilidades. Podríamos comparar la probabilidad de **si** y de **no** en un mismo estrato. Eso nos lleva a la definición de **odds** o **chance**. El odds de un suceso es

$$odds = \frac{\textit{probabilidad}}{1 - \textit{probabilidad}}$$

y toma cualquier valor mayor o igual a 0.

En el ejemplo, tenemos que para el grupo tratado el odds estimado resulta

$$0,0094/(1 - 0,0094) = 0,0094/0,9906 = 0,0095,$$

mientras que para el grupo placebo el odds estimado es

$$0,0171/(1 - 0,0171) = 0,0171/0,9829 = 0,0174.$$

En el grupo que recibió placebo la chance de infarto es 0.0174 la de no tener infarto, mientras que en el grupo tratado la chance de infarto es 0.0095 la de no tener infarto. Dicho de otra manera, en el grupo placebo la chance de no tener infarto es 57.47 veces la de infarto, mientras que en el grupo tratado, la chance de no tener infarto es 105.26 veces la de infarto.

Podríamos comparar los dos odds, por ejemplo considerando su cociente, esto da origen a

$$\begin{aligned} \theta = \text{odds ratio} &= \frac{P(Y = 1|X = 1)/P(Y = 2|X = 1)}{P(Y = 1|X = 2)/P(Y = 2|X = 2)} \\ &= \frac{\left[\frac{\pi_{11}}{\pi_{1+}} \right] / \left[\frac{\pi_{12}}{\pi_{1+}} \right]}{\left[\frac{\pi_{21}}{\pi_{2+}} \right] / \left[\frac{\pi_{22}}{\pi_{2+}} \right]} \\ &= \frac{\pi_{11}\pi_{22}}{\pi_{12}\pi_{21}} \end{aligned}$$

Como antes observamos que al ser esta medida función de $P(Y|X)$, la inferencia es válida para los tres muestreos vistos.

El EMV es

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

Las propiedades de $\hat{\theta}$ son fáciles deducir bajo muestreo multinomial, pero también son válidas con muestreo Poisson o Producto Multinomial en el que los totales marginales por filas o bien por columnas están fijos.

Como con el riesgo relativo podemos deducir un intervalo de nivel asintótico $1 - \alpha$ para $\log(\hat{\theta})$

$$\log(\hat{\theta}) \pm z_{\alpha/2} \sqrt{\widehat{V}(\log \hat{\theta})}$$

donde

$$\widehat{V}(\log(\hat{\theta})) = \frac{1}{n_{11}} + \frac{1}{n_{12}} + \frac{1}{n_{21}} + \frac{1}{n_{22}}$$

Notemos además que si intercambiamos los roles de X e Y , obtenemos

$$\theta = \frac{\pi_{11}\pi_{22}}{\pi_{12}\pi_{21}}$$

por lo que también puede ser visto como función de $P(X|Y)$, que correspondería a tener n_{+j} fijos. El hecho de que los roles de X e Y puedan ser

intercambiados es una propiedad interesante, pues puede ser de gran utilidad pues permite usar estudios restropectivos.