

SUGGERIMENTI PER L'IMPIEGO DEL TEST CHI-QUADRATO *

di Michele Zenga **

1. Introduzione e sommario

Una delle critiche più frequentemente fatte al test chi-quadrato è che per n (numerosità campionaria) grande, differenze dalla ipotesi nulla anche di lieve entità portano quasi sempre all'accettazione della ipotesi alternativa (il test sarebbe troppo sensibile).

In questa nota, dopo aver mostrato l'opportunità di misurare la distanza delle probabilità p_i ($i = 1, 2, \dots, k$) di una alternativa dalle probabilità π_i della ipotesi nulla, facendo riferimento alla funzione

$$D = \max \left\{ \frac{|p_i - \pi_i|}{\pi_i}; i = 1, 2, \dots, k \right\}, \text{ si è proceduto alla determina-}$$

zione delle numerosità campionarie necessarie per avere valori piccoli di α (probabilità errore di prima specie) e di $\max \beta$ (β è la probabilità di un errore di seconda specie; $\max \beta$ è riferito alle alternative con $D \geq \delta$, essendo $\delta > 0$).

Le numerosità campionarie così ottenute sono molto elevate ed indicano che il test chi-quadrato è meno sensibile di quanto si ritiene e che quindi le critiche di cui sopra non sembrano molto fondate.

Da ultimo, si è mostrato come utilizzare la statistica X^2 di Karl-Pearson e la funzione D per ottenere un test con errori proporzionali,

cioè un test che soddisfa la condizione $\frac{\max \beta}{\alpha} = \nu$, essendo ν un valore positivo prefissato.

2. Distanza dell'alternativa

Si abbia una multinomiale a k classi con probabilità p_1, p_2, \dots, p_k e con numero di prove n .

Si supponga di voler verificare l'ipotesi $H_0 : (p_1 = \pi_1)$ e $(p_2 = \pi_2)$

* Lavoro presentato alla XXIX riunione scientifica della Società Italiana di Statistica, Bologna 20-22 marzo 1978.

** Libera Università degli Studi di Trento - Facoltà di Economia e Commercio.

e . . . e ($p_k = \pi_k$), contro l'alternativa $H_1 : (p_1 \neq \pi_1) \text{ o } (p_2 \neq \pi_2) \text{ o } \dots$
o ($p_k \neq \pi_k$).

Come è noto, tale verifica si esegue ricorrendo alla statistica

$$X^2 = \sum \frac{(n_i - n\pi_i)^2}{n\pi_i}, \text{ essendo } n_i \text{ il numero delle volte in cui nelle } n \text{ prove}$$

si è verificata la classe i^{ma} , ovviamente $\sum n_i = n$. Sotto H_0 la statistica X^2 tende a distribuirsi, per n grande, secondo una v.c. chi-quadrato centrale con $k - 1$ gradi di libertà. Sotto una qualsiasi alternativa (p_1, p_2, \dots, p_k) la variabile X^2 tende a distribuirsi, per n grande, secondo una v.c. chi-quadrato non centrale con $k - 1$ gradi di libertà e parametro

$$\text{di non centralità } \lambda^2 = n \sum \frac{(p_i - \pi_i)^2}{\pi_i}.$$

Alcuni studiosi sostengono [1] che per valori elevati di n , anche ai più bassi livelli di significatività, il test basato su X^2 porta a respingere l'ipotesi H_0 anche per alternative assai prossime ad H_0 . Queste affermazioni sono basate, più che sull'effettivo calcolo della funzione di potenza, sulla proprietà della consistenza del test X^2 . La consistenza è però una proprietà asintotica per cui non è possibile sapere da quali valori di n si possa ritenere trascurabile la probabilità β .

Si mostrerà in seguito che, se fra l'alternativa ed H_0 vi è una piccola distanza e si vogliono avere valori di α e di β piccoli, è necessario impiegare valori molto elevati di n .

Non potendo però disporre nella pratica di valori tanto elevati di n , si impone l'esigenza di considerare le alternative aventi almeno una certa distanza da H_0 e di trascurare le alternative più vicine. Trascurare le alternative molto prossime ad H_0 non sembra possa portare a conseguenze pratiche, perché al ricercatore interessa scoprire se l'alternativa si discosti o meno di una certa entità da H_0 .

Inoltre, si tenga presente che una certa distanza fra H_0 e l'alternativa è necessaria anche per motivi strettamente tecnici, in quanto anche per valori elevati di n , se fra una alternativa ed H_0 non vi è una certa distanza, le distribuzioni di X^2 sotto le due ipotesi tendono a confondersi ed è allora difficile discriminare fra le stesse.

Vi sono diversi modi per valutare la distanza di una alternativa da H_0 . Solitamente [2] si ricorre al parametro di non centralità

$$\lambda^2 = n \sum \frac{(p_i - \pi_i)^2}{\pi_i}. \text{ Ciò è forse dovuto al fatto che la distribuzione di}$$

X^2 dipende solo da $(k - 1)$ e da λ^2 .

A fronte di questi vantaggi, λ^2 presenta alcuni inconvenienti. Innanzitutto λ^2 fa dipendere la distanza di una alternativa da H_0 dalla numerosità campionaria n . Questo inconveniente è facilmente elimina-

bile riferendosi ad $\eta^2 = \sum \frac{(p_i - \pi_i)^2}{\pi_i}$. Ma anche η^2 presenta ancora degli

inconvenienti in quanto il ricercatore non riesce a valutare che cosa voglia significare un certo valore di η^2 , essendo lo stesso, generalmente, più abituato a ragionare in termini di differenze percentuali fra p_i e π_i , piuttosto che sulle differenze al quadrato fra p_i e π_i rapportate a π_i .

Sembra, pertanto, più utile fornire un procedimento, per misurare

la distanza, basato sulle differenze percentuali assolute $\frac{|p_i - \pi_i|}{\pi_i}$.

Ad esempio, si potrebbe considerare la seguente funzione degli scarti percentuali assoluti: $D = \max \left\{ \frac{|p_i - \pi_i|}{\pi_i}; i = 1, 2, \dots, k \right\}$.

Al ricercatore, potrebbero interessare solo quelle alternative per le quali $D \geq \delta$, essendo $\delta > 0$. D'ora in avanti si indicherà con Ω_δ il sottospazio dello spazio parametrico caratterizzato da $D \geq \delta$. Affermare che $D \geq \delta$ è equivalente ad affermare che almeno una differenza percentuale assoluta

$\frac{|p_i - \pi_i|}{\pi_i}$ sia maggiore od uguale a δ .

Ovviamente, è possibile far riferimento anche ad altre funzioni delle differenze percentuali assolute per misurare la distanza fra una alternativa ed H_0 . Ad esempio, si potrebbe considerare la media aritmetica

delle percentuali assolute $\frac{1}{k} \sum \frac{|p_i - \pi_i|}{\pi_i}$ e considerare le alternative

per le quali detta media sia maggiore o uguale a v , essendo $v > 0$. Come anche, si potrebbe far riferimento ad una funzione delle differenze assolute. In questo lavoro verrà considerata la funzione:

$$D = \max \left\{ \frac{|p_i - \pi_i|}{\pi_i}; 1, 2, \dots, k \right\}.$$

Sembra ora utile avere qualche idea sulla forma della porzione dello spazio parametrico caratterizzato da $D \leq \delta$.

Lo spazio parametrico Ω può considerarsi un sottospazio dello spazio euclideo K -dimensionale E^k i cui punti (p_1, p_2, \dots, p_k) sono tali che: $\sum p_i = 1$, $p_i \geq 0$ ($i = 1, 2, \dots, k$). Sia (π_1, \dots, π_k) un punto prefissato di Ω e si definisca la funzione.

$$D = \max \left\{ \frac{|p_i - \pi_i|}{\pi_i}; i = 1, 2, \dots, k \right\}, \quad (1)$$

essendo (p_1, p_2, \dots, p_k) un punto arbitrario di Ω .

Fissato un numero reale $\delta > 0$ consideriamo il luogo dei punti che soddisfano la condizione: $\max \left\{ \frac{|p_i - \pi_i|}{\pi_i}; i = 1, 2, \dots, k \right\} = \delta$. (2)

Si controlla subito che trattasi di una ipersuperficie di dimensione $(k-1)$, e precisamente di un poliedro le cui facce sono di dimensione $(k-2)$ e sono caratterizzate dall'avere una coordinata p_i ($i = 1, 2, \dots, k$) costante. In particolare le facce sono a due a due parallele, poiché giacciono su iperpiani di equazioni: $p_i = \pi_i + \delta\pi_i$ e $p_i = \pi_i - \delta\pi_i$.

Non sempre però le facce sono $2k$, come si riscontrerà negli esempi che seguono. Tuttavia, come verrà mostrato in seguito, ciò non ha nessuna conseguenza operativa nell'impiego del test. Per esemplificare consideriamo il caso $k = 3$ con $\pi_1 = 0,4$, $\pi_2 = 0,3$ e $\pi_3 = 0,3$ e con $\delta = 0,2$.

I punti dello spazio parametrico in cui $D \leq \delta$ sono caratterizzati dal contemporaneo realizzarsi delle seguenti relazioni:

$$\begin{aligned} 0,32 &\leq p_1 \leq 0,48 \\ 0,24 &\leq p_2 \leq 0,36 \\ 0,64 &\leq p_1 + p_2 \leq 0,76 \end{aligned}$$

Conseguentemente, il contemporaneo verificarsi delle relazioni precedenti dà luogo in Ω all'esagono indicato nel grafico 1.

Non sempre però, per $k = 3$, il luogo dei punti che si ottiene ponendo $D = \delta$ dà luogo ad un esagono.

Per mostrare ciò e senza perdere in generalità si supponga $\pi_1 \geq \pi_2 \geq \pi_3$, il che implica $\pi_1 \geq 1/3$. Si ricava facilmente che se $1/3 \leq \pi_1 < 1/2$ si ha un esagono e se $\pi_1 \geq 1/2$ si ha un quadrilatero, come è indicato nel grafico 2 in cui si è supposto $\pi_1 = 0,6$, $\pi_2 = 0,2$ e $\pi_3 = 0,2$ e $\delta = 0,20$.

3. *Riformulazione del problema ed alternativa più sfavorevole.*

È possibile riformulare il problema di verifica di ipotesi riguardante la multinomiale ipotizzando che il sottospazio parametrico dell'ipotesi alternativa H_1 sia Ω_δ .

Per il problema di verifica di ipotesi di cui sopra è possibile far ancora riferimento alla statistica X^2 .

Si supponga che il valore critico della variabile test sia c . Per n sufficientemente elevato il valore di α è approssimato bene dalla probabilità che la variabile chi-quadrato con $k-1$ gradi di libertà, che si indica con $X^2_{(k-1)}$, sia maggiore o uguale a c . Analogamente per la alternativa $(p_1, p_2, \dots, p_k) \in \Omega_\delta$ e sempre per n grande il valore di β

