

Qualche riga in piú sul F-Test.

Consideriamo una serie di K esperimenti. In ognuno dei quali sono state fatte un certo numero (m_k) di osservazioni di una variabile aleatoria X ($x_{k,i}$). Ogni esperimento permette di ottenere un valor medio (\bar{X}_k) e la stima campionaria della varianza (s_k^2) per il campione dato. L'ipotesi che si vuole verificare é se i singoli esperimenti possono essere considerati come campionamenti della stessa "popolazione", quindi se i valori medi e le varianze trovate possono essere considerati come stima del valore aspettato e della varianza della popolazione globale.

	Exp ₁	Exp ₂	. . .	Exp _K
	m_1	m_2		m_K
	$x_{1,1}$	$x_{2,1}$		$x_{K,1}$
	$x_{1,2}$	$x_{2,2}$		$x_{K,2}$
	$x_{1,3}$	$x_{2,3}$		$x_{K,3}$
		
	x_{1,m_1}	x_{2,m_2}		x_{K,m_K}
medie:	\bar{X}_1	\bar{X}_2		\bar{X}_K
dev.st:	s_1	s_2		s_K

Dato un insieme di K esperimenti, ognuno numerositá m_k , assumiamo che tutte le misure vengono dalla stessa popolazione. In tal caso il valore medio, calcolato su tutte le osservazioni ($x_{j,k}$) é:

$$\bar{X} = \frac{\sum_{i=1}^{m_1} x_i + \sum_{i=1}^{m_2} x_i + \dots + \sum_{i=1}^{m_K} x_i}{m_1 + m_2 + \dots + m_K} = \frac{\sum_{k=1}^K m_k \bar{X}_k}{\sum_{k=1}^K m_k} = \frac{\sum_{k=1}^K m_k \bar{X}_k}{m}$$

con $m = m_1 + m_2 + \dots + m_K$. In pratica é una media pesata: il valor medio dell'esperimento k -esimo é pesato per la sua molteplicitá m_k . Nota: la media pesata é eguale alla media su tutte le osservazioni di tutti i campioni. Al contrario la media delle medie non é uguale alla media su tutte le osservazioni.

Prima di andare avanti notiamo che la varianza campionaria su un insieme di osservazioni é:

$$\begin{aligned} s^2 &= \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{X})^2 = \\ &= \frac{1}{m-1} \sum_{i=1}^m (x_i^2 + \bar{X}^2 - 2x_i \bar{X}) = \frac{1}{m-1} \sum_{i=1}^m (x_i^2 - x_i \bar{X}) = \\ &= \frac{1}{m-1} \left(\sum_{i=1}^m x_i^2 - m \bar{X} \frac{1}{m} \sum_{i=1}^m x_i \right) = \frac{1}{m-1} \left(\sum_{i=1}^m x_i^2 - m \bar{X}^2 \right) = \\ &= \frac{m}{m-1} \left(\frac{1}{m} \sum_{i=1}^m x_i^2 - \bar{X}^2 \right) = \\ &= \frac{m}{m-1} (\bar{X}^2 - \bar{X}^2) \end{aligned}$$

Quindi la varianza campionaria può essere calcolata come somma del quadrato degli scarti diviso $m-1$, oppure come differenza tra la media dei quadrati delle osservazioni e il quadrato della media, normalizzata per $m/m-1$. Questa espressione può essere utile per semplificare le espressioni: ad esempio, avendo a che fare con K esperimenti ognuno con un valor medio \bar{X}_k , la migliore stima della media è \bar{X} , come visto sopra (e non la media delle medie, infatti è intuitivo che campioni con una numerosità maggiore forniscano una migliore stima del valore aspettato vero della popolazione), l'errore standard sulla media si può stimare come:

$$\begin{aligned} s_{\bar{X}}^2 &= \frac{1}{K-1} \sum_{k=1}^K (\bar{X}_k - \bar{X})^2 = \frac{1}{K-1} \left(\sum_{k=1}^K (\bar{X}_k^2 + \bar{X}^2 - 2\bar{X}\bar{X}_k) \right) = \\ &= \frac{1}{K-1} \left(\sum_{k=1}^K \bar{X}_k^2 + K\bar{X}^2 - 2\bar{X} \sum_{k=1}^K \bar{X}_k \right) = \\ &= \frac{1}{K-1} \left(\sum_{k=1}^K \bar{X}_k^2 + K\bar{X}^2 - 2K\bar{X}^2 \right) = \\ &= \frac{1}{K-1} \left(\sum_{k=1}^K \bar{X}_k^2 - K\bar{X}^2 \right) = \\ &= \frac{K}{K-1} (\bar{\bar{X}}^2 - \bar{X}^2) \end{aligned}$$

dove $\bar{\bar{X}}^2$ è la media dei quadrati delle medie campionarie di ogni esperimento, invece \bar{X}^2 è il quadrato della media stimata utilizzando tutti i dati di tutti gli esperimenti.

Se si può vuole stimare la varianza di tutte le osservazioni possiamo utilizzare due approcci: nel primo caso (s_{tra}^2) assumiamo che, nell'ipotesi da verificare che tutti i dati vengano dalla stessa popolazione, tutti gli m_k valori nell'esperimento k -esimo abbiano uno stesso valore medio atteso, la cui miglior stima è appunto \bar{X}_k . Quindi stimiamo s_{tra}^2 come media dei quadrati degli scarti tra media dell'esperimento k -esimo e \bar{X}_k . Dal momento che esperimenti con maggiore numerosità sono probabilmente più precisi, pesiamo i quadrati degli scarti per la numerosità del campione k -esimo:

$$\begin{aligned} s_{tra}^2 &= \frac{1}{K-1} \sum_{k=1}^K m_k (\bar{X}_k - \bar{X})^2 = \frac{1}{K-1} \left(\sum_{k=1}^K m_k (\bar{X}_k^2 + \bar{X}^2 - 2\bar{X}\bar{X}_k) \right) = \\ &= \frac{1}{K-1} \left(\sum_{k=1}^K m_k \bar{X}_k^2 + \bar{X}^2 \sum_{k=1}^K m_k - 2\bar{X} \sum_{k=1}^K m_k \bar{X}_k \right) = \\ &= \frac{1}{K-1} \left(\sum_{k=1}^K m_k \bar{X}_k^2 + m\bar{X}^2 - 2m\bar{X}^2 \right) = \\ &= \frac{1}{K-1} \left(\sum_{k=1}^K m_k \bar{X}_k^2 - m\bar{X}^2 \right) \end{aligned}$$

Notiamo che, se le numerosità dei campioni sono le stesse ($m_1 = m_2 = \dots m_k = m_o$), si ha:

$$s_{tra}^2 = \frac{m_o}{K-1} \sum_{k=1}^K (\bar{X}_k - \bar{X})^2 = \frac{1}{K-1} \left(m_o \sum_{k=1}^K \bar{X}_k^2 - m_o K \bar{X}^2 \right)$$

$$= m_o \frac{K}{K-1} \left(\frac{1}{K} \sum_{k=1}^K \bar{X}_k^2 - \bar{X}^2 \right) = m_o \frac{K}{K-1} (\overline{X^2} - \bar{X}^2)$$

dove $\overline{X^2}$ é la media dei quadrati delle medie campionarie di ogni singolo esperimento mentre \bar{X}^2 é la media stimata utilizzando tutti i dati di tutti gli esperimenti. Per campioni della stessa numerosit  si ha: $s_{tra}^2 = m_o s_{\bar{x}}^2$ che giustifica la relazione tra stima della varianza e errore standard della media: $s_{\bar{x}}^2 = s^2 / \sqrt{m_o}$

Un secondo modo per stimare la varianza delle osservazioni é come media pesata delle varianze campionarie, in cui il peso é il numero di gradi di libert  per ognuna di esse:

$$s_{intra}^2 = \frac{\sum_{k=1}^K (m_k - 1) s_k^2}{\sum_{k=1}^K (m_k - 1)} = \frac{\sum_{k=1}^K (m_k - 1) s_k^2}{m - K}$$

che, per campioni della stessa numerosit , diventa:

$$s_{intra}^2 = \frac{1}{K} \sum_{k=1}^K s_k^2$$

A volte s_{intra}^2 si trova anche scritta con un formalismo un p  pi  complesso:

$$\begin{aligned} s_{intra}^2 &= \frac{\sum_{k=1}^K (m_k - 1) s_k^2}{\sum_{k=1}^K (m_k - 1)} = \frac{1}{m - K} \sum_{k=1}^K (m_k - 1) \frac{m_k}{m_k - 1} (\overline{X_k^2} - \bar{X}_k^2) = \\ &= \frac{1}{m - K} \sum_{k=1}^K m_k \left(\frac{1}{m_k} \sum_{i=1}^{m_k} x_{k,i}^2 - \bar{X}_k^2 \right) = \\ &= \frac{\sum_{k=1}^K \sum_{i=1}^{m_k} x_{k,i}^2 - \sum_{k=1}^K m_k \bar{X}_k^2}{m - K} \end{aligned}$$

Se i campioni vengono effettivamente dalla stessa popolazione (ipotesi da verificare), allora il rapporto: $F = s_{tra}^2 / s_{intra}^2 \sim 1$. In generale si pu  dimostrare che $s_{intra}^2 < s_{tra}^2$ e che ci  é tanto pi  vero quanto pi  sono diverse le popolazioni d'origine dei campioni.

La funzione F é una variabile aleatoria che segue una distribuzione nota come distribuzione F, la cui forma dipende dai gradi di libert  al numeratore ($\nu_1 = K - 1$) e al denominatore ($\nu_2 = m - K$). Questa funzione si trova tabulata in funzione di ν_1 , ν_2 e del livello di confidenza: $F(\nu_1, \nu_2, \alpha)$. Il valore riportato nelle tabelle indica il valore critico F_c con il seguente significato: nell'ipotesi che i campioni siano tutti provenienti da una stessa popolazione la probabilit  di osservare valori $F > F_c$ é minore di α mentre la probabilit  di osservare valori di $F < F_c$ é maggiore di $1 - \alpha$. Quindi se si ottiene un valore $F < F_c$ possiamo dire che l'ipotesi é verificata con un margine di errore (rischio) pari ad α (o che é la stessa cosa l'ipotesi é verificata con probabilit  $1 - \alpha$). Quindi posso assumere che le differenze eventualmente osservate nei diversi campioni non sono significative da un punto di vista statistico.

Al contrario se si ottiene un valore $F > F_c$ possiamo dire che l'ipotesi non é verificata, quindi i campioni presentano diversit  statisticamente significative, con un margine di errore (rischio) pari ad α (o che é la stessa cosa l'ipotesi non é verificata con una confidenza $1 - \alpha$).

In Excel la funzione `INV.F(n1,n2,alpha)` restituisce il valore F_c per un livello di confidenza α .