

ANOVA: Analisi della varianza (ANalysis Of VAriance)

generalizzazione del test del confronto tra due medie al caso di g popolazione

→ studio degli effetti di un fattore sperimentale (variabile qualitative) su una variabile di risposta (variabile quantitative)

MEMO: studio delle dipendenze in media

GRUPPI	CARDINALITÀ	MEDIE	DEVIANZE	VARIANZE
1	m_1	\bar{x}_1	Dev_1	s_1^2
2	m_2	\bar{x}_2	Dev_2	s_2^2
⋮	⋮	⋮	⋮	⋮
i	m_i	$\bar{x}_i = \sum_{j=1}^{m_i} \frac{x_{ij}}{m_i}$	$Dev_i = \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2$	$s_i^2 = \frac{Dev_i}{m_i - 1}$
⋮	⋮	⋮	⋮	⋮
G	m_G	\bar{x}_G	Dev_G	s_G^2

(M)

$$\sigma^2 = \sigma_{EST}^2 + \sigma_{INT}^2 \quad \text{proprietà di decomposizione della varianza}$$

$$Dev(X) = Dev_{EST} + Dev_{INT} \quad \text{in termini di devianze}$$

EST \rightarrow esterna o tra i gruppi

INT \rightarrow interna o all'interno dei gruppi

$$\text{Der}(X) = \sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2$$

varianza complessiva del carattere

$$\text{Der}_{\text{EST}} = \sum_{i=1}^g (\bar{x}_i - \bar{x})^2 n_i$$

misura della differenza tra i gruppi (tra le medie di gruppo e la media complessiva)

$$\text{Der}_{\text{INT}} = \sum_{i=1}^g \sum_{j=1}^{m_i} (x_{ij} - \bar{x}_i)^2 = \sum_{i=1}^g \text{Der}_i = \sum_{i=1}^g s_i^2 (n_i - 1)$$

misura delle varianze interne ai gruppi

Le ipotesi dell'ANOVA:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_g \quad \text{le } g \text{ popolazioni hanno le stesse medie}$$

$$H_1: \text{almeno una } \mu_i \text{ è diversa dalle altre}$$

NOTA TECNICA: l'ANOVA assume che le g popolazioni siano normali e che abbiano la stessa varianza
(ipotesi di omoschedasticità)

i due estimatori di σ^2 sono ottenuti sfruttando Dev_{NT} e Dev_{EST}

$$\sigma_1^2 = \sigma_2^2 = \dots = \sigma_g^2 = \sigma^2$$

la verifica di H_0 consiste nel confrontare due estimatori della varianza comune, uno valido sia sotto H_0 che sotto H_1 e l'altro valido solo sotto H_0

- STIMATORE 1 PER σ^2 (valido sia sotto H_0 che sotto H_A)



ciascuna delle varianze campionarie s_i^2 è uno stimatore corretto della varianza della popolazione



a partire dai G stimatori disponibili $s_1^2, s_2^2, \dots, s_G^2$
si ottiene un unico stimatore che m'è le medie perate:

$$\frac{\sum_{i=1}^G s_i^2 (m_i - 1)}{m - G} = \frac{\sum_{i=1}^G \text{Der}_i}{m - G} = \frac{\text{Der}_{INT}}{m - G}$$



questo stimatore si ottiene senza sapere nulla sulla verità o sulla falsità di H_0

- STIMATORE $\hat{\sigma}$ PER σ (valido sotto H_0)



sotto H_0 si ha che $\mu_i = \mu \quad \forall i = 1, g$



le g medie campionarie $\bar{X}_1, \bar{X}_2, \dots, \bar{X}_g$ si distribuiscono come una normale con media μ e varianza $\frac{\sigma^2}{m_i}$



la varianza campionaria di questi dati:

$$\bar{s}^2 = \frac{\sum_{i=1}^g (\bar{X}_i - \bar{X})^2 m_i}{g-1} = \frac{\text{Dev est}}{g-1}$$

è uno stimatore di σ^2

Da cui si ottiene la statistica test F :

$$F = \frac{\text{Dev}_{\text{TRA}} / (g - 1)}{\text{Dev}_{\text{ENTRO}} / (m - g)} \sim F_{g-1, m-g}$$

stima corretta di σ^2 sotto H_0 (tende a sovrastimare
se H_0 è falsa)

stima corretta di σ^2 sia sotto H_0 che sotto H_1

il rapporto tende a essere grande quando H_0 è falsa

(del resto se la Dev_{TRA} è molto alta rispetto alle Dev_{ENTRO} è un
segnale che le medie dei gruppi sono molto diverse rispetto alle
medie complessive)

Le statistiche dell'ANOVA sono solitamente riassunte in una tavola di questo tipo:

FONTE DI VARIABILITÀ	DEVIANZE (SOMME DEI QUADRATI)	g.d.l.	VARIANZE (MEDIIE DEI QUADRATI)	F	p-value
tra i gruppi (spiegata)	Dev _{EST}	G - 1	Dev _{EST} / (G - 1)	Dev _{EST} / (G - 1) / Dev _{INT} / (m - G)	p-value associato alla statistica F
entro i gruppi (residua)	Dev _{INT}	m - G	Dev _{INT} / (m - G)		
totale	Dev _{TOT}	m - 1	Dev _{TOT} / (m - 1)		

La regola di decisione dell'ANOVA è sempre di tipo univorzionale α :

$$\frac{A}{F_{G-1, m-G; \alpha}}$$

si può procedere confrontando il valore della statistica F con le soglie o confrontando il p-value con l' α

ESERCIZIO (ANOVA)

In una casa editrice ci sono tre redazioni (A, B e C). A partire dalle 3 redazioni sono stati selezionati casualmente 4, 5 e 6 redattori.

L'esperienza in mesi dei redattori è riportata nella seguente tabella:

REDAZIONE	ESPERIENZA
A	25, 27, 29, 31
B	23, 23, 24, 26, 29
C	25, 26, 27, 29, 34, 36

Costruire la tavola ANOVA e, assumendo che le tre popolazioni si distribuiscono normalmente, sottoporre a verifica l'ipotesi nulla che non vi sia differenza tra le esperienze medie nelle tre redazioni al livello di significatività $\alpha=0.05$.

Statistiche di base

REDAZIONE A	REDAZIONE B	REDAZIONE C
$m_A = 4$	$m_B = 5$	$m_C = 6$
$\bar{x}_A = 28$	$\bar{x}_B = 25$	$\bar{x}_C = 23,5$
$Dev_A = 20$	$Dev_B = 26$	$Dev_C = 101,5$

$$\bar{x} = \frac{28 \times 4 + 25 \times 5 + 23,5 \times 6}{15} = 27,6$$

$$Dev_{TRA} = (28 - 27,6)^2 \times 4 + (25 - 27,6)^2 \times 5 + (23,5 - 27,6)^2 \times 6 = 56,1$$

$$Dev_{ENTRO} = 20 + 26 + 101,5 = 147,5$$

$$Dev_{TOT} = Dev_{TRA} + Dev_{ENTRO} = 56,1 + 147,5 = 203,6$$

Fonte di variazione	DEV (SQ)	golf	VAR (MQ)	F
Tra i gruppi	56,1	2	28,05	2,28
Entro i gruppi	147,5	12	12,29	
TOTALE	203,6	14		

I passi del test ANOVA

$$\textcircled{1} \quad H_0: \mu_A = \mu_B = \mu_C$$

$H_1:$ almeno una μ_i è diversa dalle altre

$$\textcircled{2} \quad \alpha = 0,05$$

$$\textcircled{3} \quad F = \frac{MQ_{TRA}}{MQ_{ENTRO}} \sim F_{q-1, n-q}$$

\textcircled{4}

$$A + R$$

$$F_{2, 12; 0,05}$$

$$3,89$$

\textcircled{5}

$$F \in A$$



$$2,28 < 3,89$$

non si rifiuta H_0