

Statistica della Formazione

Slides

A.A. 2020-2021

Docente: ANNA LINA SARRA

Modulo 4 : I modelli di risposta all'item



Valutazione quantitativa degli apprendimenti

La valutazione quantitativa degli apprendimenti trova il suo fondamento nell'approccio teorico del *modello funzionalista*.

Quando parliamo del **modello teorico di riferimento**, consideriamo una particolare modalità di approccio alla valutazione costituita da una propria filosofia di pensiero, da uno substrato teorico specifico, da principi e concetti a cui richiamarsi, in modo esplicito o implicito, per affrontare empiricamente la valutazione.

Il modello funzionalista

Il modello funzionalista, nato dall'esigenza di abolire azioni didattiche basate sulla soggettività e sull'ambiguità, si caratterizza per un approccio alla valutazione di tipo quantitativo.

Si fonda, quindi, sulla verifica della congruità tra obiettivi programmati e risultati conseguiti.

Nella valutazione quantitativa, infatti, l'interesse è centrato sulla verifica della conformità tra obiettivi e risultati, e non sull'analisi del processo che sottende a questa relazione.

.

La valutazione quantitativa deve garantire



L'affidabilità di uno strumento

Quando si ripete più volte una misura e si ottiene lo stesso risultato si dice che la misura è affidabile, attendibile o fedele.

L'affidabilità si riferisce quindi alla costanza della misura di una data prestazione. L'importanza del concetto di affidabilità deriva dal fatto che, se uno strumento è attendibile nei risultati, si può essere sicuri che le variazioni che si verificano nei dati raccolti non dipendono da imperfezioni dello strumento utilizzato, ma dal mutare del fenomeno (ERRORE CASUALE).

Aldo Visalberghi a tal proposito scrive: «Non c'è nulla di assolutamente misurabile. Se molte misurazioni fisiche risultano perfettamente uguali ciò non significa che la nostra capacità misurativa al riguardo è perfetta, bensì, tutto al contrario, che il nostro strumento di misura non è abbastanza sensibile per il lavoro che stiamo facendo. Ogni misurazione è una media» (1965, p.77) .

Aldo Visalberghi (Trieste, 1° agosto 1919 – Roma, 12 febbraio 2007)
è stato un [pedagogista](#), [accademico](#), [politico](#), [partigiano](#) e [antifascista](#) italiano.

Item-analisi

Chiamiamo *item-analisi* l'insieme delle **tecniche che permettono di ricavare informazioni sulla affidabilità di una prova nel suo complesso** e sul funzionamento di **ciascuna delle domande proposte**.

Le tecniche di item-analisi muovono tutte dall'assunto che, se la prova nel suo complesso costituisce una misura di una dimensione unitaria, le singole domande e i singoli soggetti dovranno avere un comportamento coerente.

Che cos'è l'analisi degli item (Item analysis) in generale?

- ✓ **L'analisi degli item** fornisce
- ✓ un modo per **misurare la qualità delle domande** - vedere quanto siano appropriate per gli intervistati e quanto bene misurino la loro abilità / caratteristica.
- ✓ un modo per riutilizzare gli items più e più volte in diversi test con una conoscenza preliminare delle loro prestazioni, creando una popolazione di domande con proprietà note (item bank)

Tratti latenti (latent trait)

LATENT TRAIT –Lord et al. (1968)

Costrutti non direttamente osservabili ma desumibili da un insieme di osservazioni (prove attitudinali, questionari, prove fisiche, mentali, ecc.) che consentono appunto di effettuare una stima dei fenomeni stessi.

Tratti latenti (latent trait)

Nel linguaggio di uso comune noi riusciamo facilmente a descrivere questo costrutto elencandone consapevolmente gli attributi specifici **ma non possediamo uno strumento di misura che sia in grado di stimare direttamente la loro consistenza**, in quanto la variabile latente esprime un concetto astratto e non una dimensione fisica concreta.

ABILITA':

Un costrutto, ossia un insieme di concetti astratti che indicano un aspetto della vita intellettuale del soggetto, non osservabile direttamente ma inferito a partire da una serie di indicatori empirici osservabili

Operazionalizzazione

I tratti latenti sono dimensioni teoriche (costrutti) non osservabili direttamente, ma misurabili mediante l'individuazione di indicatori osservabili di cui sono l'espressione.

L'indicatore è una variabile osservabile che si ipotizza cogliere il costrutto o parte di esso.

La scelta degli indicatori non è ovvia.

C'è sempre una componente arbitraria poiché dipendono da un lato dalle teorie dei ricercatori e dall'altro dagli strumenti adottati per misurarli.

Indicatori riflettivi ed indicatori formativi

Il rapporto tra il tratto latente ed indicatori può verificarsi in due situazioni:

1. Gli indicatori riflettono il costrutto (**INDICATORI RIFLETTIVI**)
2. Gli indicatori, al contrario, formano o causano il costrutto (**INDICATORI FORMATIVI**)

Gli indicatori **riflettono** (manifestazione osservabile) il costrutto (*scala*; es., “**sensibilità al rumore**” → sono svegliato dal rumore, mi abituo al rumore, mi irrita se i vicini sono rumorosi, ecc.).

Gli indicatori **formano** il costrutto (*indice*; es., “**stress**” ← trasloco, divorzio, nascita di un figlio, lutto, ecc.).

Modulo 4.1: La teoria classica dei test (TCT)



$$\mathbf{X} = \mathbf{T} + \mathbf{E}$$

Observed Score = True Score + Error

La Teoria Classica dei Test (TCT)

La Teoria Classica dei Test (Classical Test Theory) nasce alla fine dell'Ottocento (Alfred Binet e altri, 1894) e si sviluppa con il duplice interesse

- ✓ di misurazione di un costrutto
- ✓ e validazione dello strumento utilizzato nella misurazione dello stesso

L'impiego su vasta scala e lo sviluppo della TCT ha inizio negli anni Trenta, anche se la formalizzazione dell'equazione fondamentale su cui tale teoria si basa viene proposta da Spearman qualche decennio prima.

La teoria classica dei test (TCT)

Questo modello consiste nel **sostenere che il punteggio che una persona ottiene nel test, che denominiamo punteggio empirico**, e che solitamente è indicato dalla lettera X, è formato da due componenti:

$$X = T + E$$

• IL PUNTEGGIO VERO

• L'ERRORE NON SISTEMATICO

L'errore può essere dovuto a molte cause che non possiamo controllare. Per questo la TCT si **preoccupa di determinare con precisione l'errore di misurazione**.

Errori

L'errore è parte integrante del processo di misurazione.

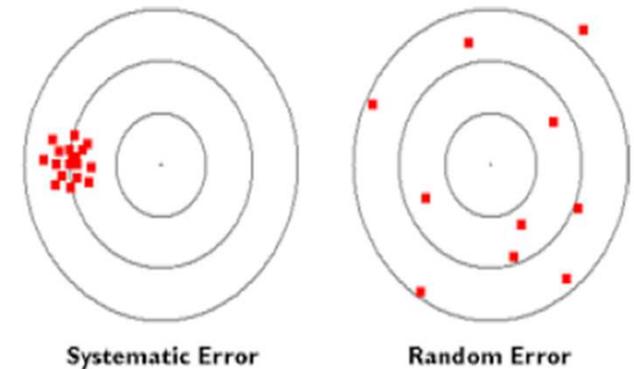
• Si possono distinguere due tipi di errori:

- ERRORI CASUALI

sono errori che variano in maniera casuale tra le diverse misurazioni. Non predicibili

- ERRORI SISTEMATICI (BIAS)

sono errori che si presentano in maniera costante e predicibile



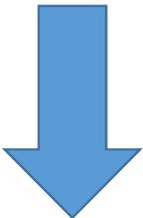
L'errore di misurazione casuale e non sistematico

L' errore di misurazione casuale e non sistematico non è una proprietà della caratteristica misurata ma è conseguente alla misurazione effettuata.

Esso è, quindi, inversamente proporzionale all'affidabilità:



MAGGIORE è LA COMPONENTE DI ERRORE

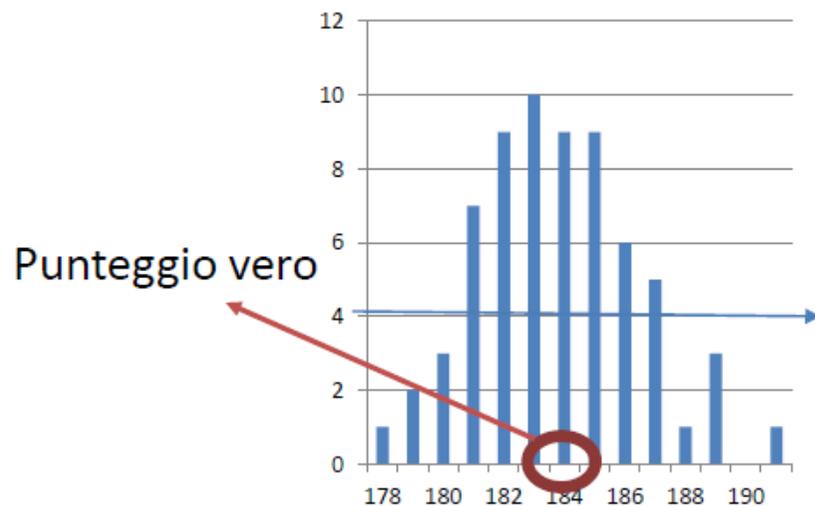


MINORE SARA' L'AFFIDABILITA' DELLA MISURAZIONE

Teoria classica dei test (TCT)

Immaginiamo che un test venga somministrato N volte (N è un valore tendenzialmente infinito) allo stesso soggetto.

Tutti i punteggi osservati possono essere quindi considerati come una variabile che assume dei valori secondo una distribuzione di probabilità che contiene il punteggio vero.

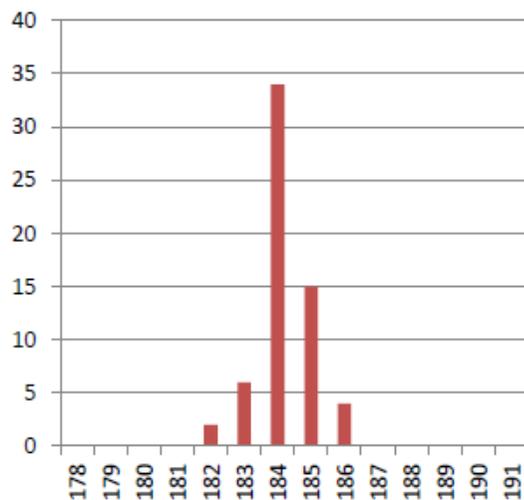


Legge dei grandi numeri: la media che calcoliamo a partire da un *numero sufficiente* di campioni sarà *sufficientemente vicina* alla media vera

Teoria classica dei test (TCT)

Di conseguenza, il punteggio di un soggetto ad un test è un campione di una popolazione di infiniti possibili punteggi affetti da errori casuali.

La popolazione contiene il valore vero e assume una forma normale.



Teorema limite centrale: la distribuzione della somma (o media) di un numero elevato di variabili casuali indipendenti e identicamente distribuite tende a distribuirsi come una distribuzione normale, indipendentemente dalla distribuzione delle singole variabili

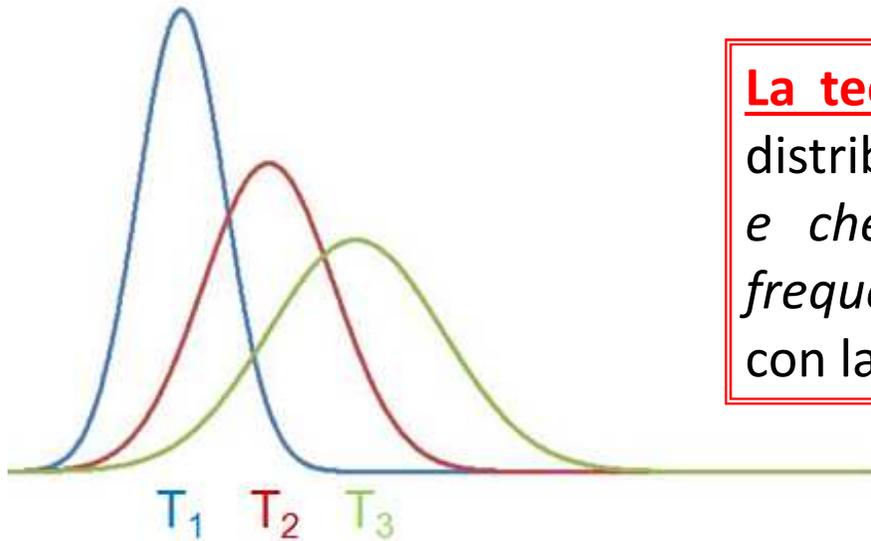
Teoria classica dei test (TCT)

Il punteggio vero rappresenta il **valore atteso** della distribuzione di probabilità associata ai punteggi osservati.

$$E(X) = T$$


OPERATORE DI VALORE ATTESO (MEDIA)

Punteggio vero



La teoria classica della misurazione assume che la distribuzione di tutti i valori così registrati sia *normale* e che il valore vero sia quello che presenta la frequenza più alta o, meglio, che tale valore sia quello con la probabilità più alta di avvicinarsi a quello vero.

Il vero punteggio di una persona verrebbe ottenuto facendogli rispondere a tutti gli item nell' "universo" degli item. Nel test vengono visualizzate **solo le risposte al campione di item.**

Ne consegue che...

Se le misure si distribuiscono secondo una curva normale allora dopo moltissime, infinite, misurazioni, il **valore medio di X corrisponde al valore vero di X**.

Da $X=T+E$ segue che T rappresenta la parte attendibile del punteggio, dato che $E(X)=T$.

In una serie infinita di misurazioni l'errore casuale scompare:

$$\lim_{n \rightarrow \infty} \left(\frac{\sum_{i=0}^n E}{n} \right) = 0 \quad E(E) = 0$$

Ne consegue che...

Dato che l'errore casuale

$$E = X - T$$

$E(E)$ avrà una distribuzione normale centrata sullo 0

Se questo non accade l'errore non è casuale

Ne consegue che...

Inoltre, dalla precedente assunzione segue che la covarianza tra punteggio vero ed errore è nulla.

$$\sigma_{TE} = 0$$

E che covarianza tra due distribuzioni di errore è nulla

$$\sigma_{EX}\sigma_{EY} = 0$$

Le tre ipotesi della TCT

- I. **Il punteggio vero (T) è l'aspettativa matematica del punteggio empirico: $T = E(X)$**
- II. ***Il valore del punteggio vero (T) è indipendente dall'errore di misurazione.***
- II. ***Gli errori di misura in un test concreto non sono legati agli errori di misura di un altro test diverso.***

Gli errori commessi in un'occasione non influenzerebbero quelli commessi in un'altra occasione.

In sintesi.....

LA TEORIA CLASSICA DEI TEST ASSUME CHE

- 1. L'ERRORE È NORMALMENTE DISTRIBUITO**
- 2. INCORRELATO CON PUNTEGGIO VERO**
- 3. HA MEDIA ZERO**

Estensione a più soggetti-CASO REALE

I punteggi differiscono tra soggetti, non soltanto per l'errore di misurazione E , ma anche e soprattutto per le **differenze individuali**, che si riflettono nelle differenze tra i punteggi osservati che non dipendono dall'errore di misurazione.

- I punteggi veri ottenibili per ogni individuo risultano uguali solo nel caso banale in cui **il costrutto oggetto di studio non ha variabilità**.

Estensione a più soggetti-CASO REALE

Somministrando N test paralleli agli stessi soggetti, si può assumere che le medie dei punteggi dei soggetti calcolate su ogni test siano uguali tra loro:

$$E(X1) = E(X2) = \dots = E(XN)$$

dove $E(X1)$, per esempio, rappresenta la media dei punteggi sugli n soggetti al primo test.

- Inoltre, queste medie coincidono con la media dei punteggi veri :

$$E(X1) = \dots = E(XN) = E(T)$$

- Come conseguenza, le medie di tutti gli errori sono uguali a zero

$$E(E1) = E(E2) \dots E(EN) = 0$$

Attendibilità

Attendibilità, affidabilità e fedeltà

sono tre sinonimi utilizzati per riferirsi, in ambito statistico, al grado di accuratezza e precisione di una procedura di misurazione.

- Si dice allora che un test è affidabile, quando si può affermare che i punteggi ottenuti da un gruppo di soggetti allo stesso

“sono coerenti, stabili nel tempo e costanti dopo molte somministrazioni e in assenza di cambiamenti evidenti quali variazioni psicologiche e fisiche degli individui che si sottopongono al test, o anche all’ambiente in cui questo ha luogo”

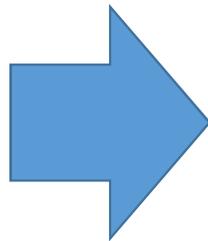
(Pedrabissi, Santinello, 1997)

Attendibilità

Abbiamo detto che ogni misurazione è affetta da errori.

Gli errori sistematici possono essere identificati perché predicibili.

Gli errori casuali non sono predicibili.



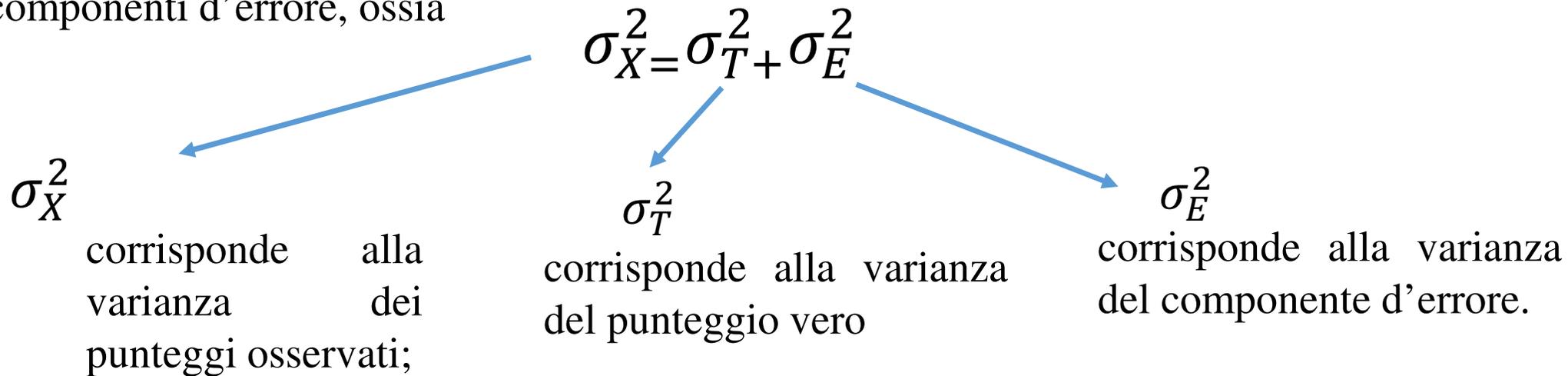
Nella equazione $X = T + E$ come possiamo separare T da E in modo da capire quanta informazione vera (T) è contenuta in X ?

L'attendibilità misura il grado di coerenza e di stabilità di un test o anche il grado di precisione con cui una scala misura un costrutto.

La varianza

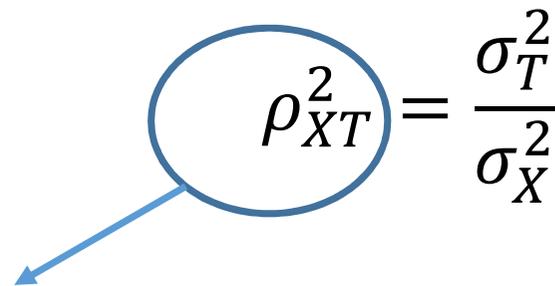
La varianza sarà dovuta a componenti vere e a componenti d'errore, cioè dipenderà dal diverso grado di possesso di una certa caratteristica da parte dei soggetti e da fattori occasionali estranei al test e non controllabili.

Possiamo allora affermare che la varianza totale della distribuzione dei punteggi equivale alla somma della varianza delle componenti vere e di quella delle componenti d'errore, ossia



L'attendibilità

L'attendibilità di un test può essere definita allora come la porzione di varianza vera rispetto alla varianza totale di una distribuzione di punteggi, cioè il **rapporto tra varianza dei punteggi veri e varianza dei punteggi osservati** può essere definito come una valutazione del livello di affidabilità,


$$\rho_{XT}^2 = \frac{\sigma_T^2}{\sigma_X^2}$$

È il coefficiente di attendibilità del test

L'analisi classica degli item

Comprende:

- Frequenze di risposta
- Indice di distrattività
- Indice di difficoltà
- Indice di discriminatività
- Correlazione item-totale
- Alfa di Cronbach

Frequenze di risposta

Quante volte ricorre una risposta in un quesito, per tutte le risposte possibili, per tutti i quesiti.

Possono essere

- ✓ ASSOLUTE
- ✓ RELATIVE
- ✓ CUMULATE
- ✓ PERCENTUALI

Quesito 1	Fq assoluta	Fq relativa	Fq cumulata	Fq %
A	6	0,188	6	19%
B	7	0,219	13	22%
C	4	0,125	17	12%
D (corretta)	15	0,469	32	47%
totale	32	1	-	100%

Indice di distrattività

L'indice di distrattività individua la capacità dei singoli distrattori di far deviare dalla risposta corretta.

La distrattività del quesito si misura tenendo conto della distribuzione delle risposte errate sui distrattori.

Nell'ipotesi di quesiti con possibilità di scelta tra quattro opzioni (una risposta esatta e tre distrattori), la situazione ideale è quella per cui tutti gli item hanno la stessa capacità di attrarre risposte e quindi la stessa frequenza.

Se uno dei distrattori è palesemente errato e nessuno lo sceglie, le opzioni vere per il rispondente rimangono ridotte a tre e il ruolo della casualità nella scelta esatta sarebbe più forte, con un abbassamento del grado di validità del quesito.

Indice di distrattività

L'**indice di distrattività** si calcola, sul complesso degli errori, quanti, per ciascun *item*, si riferiscono a ciascun distrattore.

Indice di distrattività = D/E_t

dove:

D = numero alunni che hanno scelto i diversi distrattori

E = errori totali commessi

Ad esempio si supponga che un test sia svolto da 80 studenti e il quesito a risposta multipla n. 1 abbia avuto il seguente esito:

- risposta esatta a: n. 50
- distrattore b: n. 20
- distrattore c: n. 8
- distrattore d: n. 2

Totale risposte fornite ai distrattori = n. 30

$D_b = 20/30 = 0,67 \rightarrow$ efficace

$D_c = 8/30 = 0,26 \rightarrow$ abbastanza efficace

$D_d = 2/30 = 0,07 \rightarrow$ inefficace

Indice di **distrattività**

- I. ■ Si adopera per valutare il funzionamento dei quesiti a scelta multipla
 - II. ■ Misura la capacità di ogni singolo distrattore di far deviare dalla risposta corretta a ciascun quesito
 - III. ■ È la percentuale di studenti che sceglie ciascun distrattore.
 - IV. ■ ..senza includere le risposte lasciate in bianco, le doppie risposte o quelle incomprensibili
- Se $N < 100$ NON si usa la frequenza percentuale ma direttamente le frequenze assolute

Come si interpreta la distrattività

- I. Serve ad individuare quei distrattori poco o per nulla attrattivi
- II. Teoricamente, la risposta corretta dovrebbe avere sempre la frequenza assoluta più elevata e i distrattori dovrebbero dividere equamente i restanti rispondenti
- III. Se un distrattore non viene scelto vuol dire che non è plausibile (molto spesso per ragioni indipendenti dai contenuti) e rende per questo più facile il quesito.

Dicotomizzare le risposte

Per calcolare gli' indici di DIFFICOLTÀ/FACILITÀ e di DISCRIMINATIVITÀ occorre DICOTOMIZZARE le risposte.

1=risposta esatta
0= risposta errata

Studente	D1	D2	D3	D4	D5	Studente	D1	D2	D3	D4	D5
A1	A	C	C	SI	V	A1	1	1	1	1	0
A2	B	B	C	SI	F	A2	0	0	1	1	1
A3	B	C	D	NO	V	A3	0	1	0	0	0
A4	B	A	C	NO	F	A4	0	0	1	0	1
A5	C	C	C	SI	F	A5	0	1	1	1	1
A6	B	C	D	SI	F	A6	0	1	0	1	1
A7	D	D	A	SI	V	A7	0	0	0	1	0
A8	C	A	B	SI	F	A8	0	0	0	1	1
A9	D	C	C	NO	F	A9	0	1	1	0	1
A10	A	A	C	NO	V	A10	1	0	1	0	0
CHIAVE	A	C	C	SI	F						

Difficoltà e discriminatività

Indici di base nell'analisi classica dell'andamento dei quesiti

INDICE DI DIFFICOLTÀ'

$$Df_i = \frac{nE_i}{N}$$

INDICE DI FACILITÀ

$$Fc_i = \frac{nC_i}{N}$$

$$Fc_i = 1 - Df_i$$

Esempio

oscillazione tra 0 e 1

	Quesito 1	Quesito 2	Quesito 3	Quesito 4
Corrette	20	22	24	4
Errate	5	3	1	21
Totale	25	25	25	25
Difficoltà	0,2	0,12	0,04	0,84
Difficoltà %	20%	12%	4%	84%
Item facility	80%	88%	96%	16%

Quesito più difficile

Quesito più facile

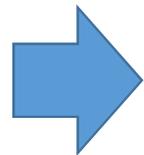
Come si interpreta l'indice di difficoltà?

Non ci sono regole sempre valide, tutto dipende dagli scopi della prova.

Solitamente se

✓ se <0.3 è troppo FACILE

✓ se >0.7 è troppo DIFFICILE



All'interno della stessa prova ci devono essere sia quesiti facili che quesiti più difficili.

Indice di difficoltà

Secondo la TCT, la scelta che viene usualmente fatta è quella di una dispersione moderata e simmetrica del livello di difficoltà attorno ad un valore leggermente superiore al valore che sta a metà tra il livello $1/n$, dove n è il numero delle alternative di risposta all'item, e il punteggio pieno.

Nel caso di un questionario a quattro alternative di risposta, il livello del caso per ogni item è pari a $1/4 = 0,25$. Il livello ottimale di difficoltà media sarà, quindi, $(0,25+1,00)/2=0,62$.

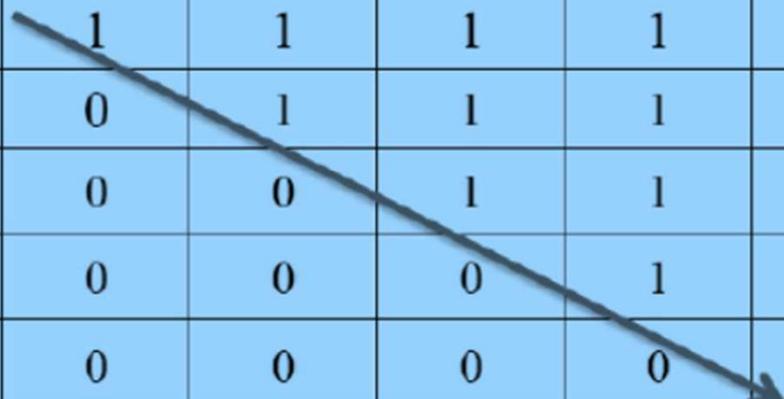
Nel caso di item con sole due alternative di risposta (es.Sì/No,Vero/Falso,ecc.) il livello ottimale di difficoltà media è $(0,50+1,00)/2=0,75$.

Aspetti critici del livello di difficoltà nella TCT

Funzionamento ideale della difficoltà

(Losito pag.107)

	Q1	Q2	Q3	Q4	Q5	%
Studente 1	1	1	1	1	1	100
Studente 2	0	1	1	1	1	80
Studente 3	0	0	1	1	1	60
Studente 4	0	0	0	1	1	40
Studente 5	0	0	0	0	1	20
Df	0,80	0,60	0,40	0,20	0	



Aspetti critici del livello di difficoltà nella TCT

Se aggiungiamo uno studente con un andamento di risposta anomalo, la situazione non è più facilmente interpretabile.

Gli studenti 4 e 6 hanno lo stesso punteggio % ma hanno risposto a quesiti di difficoltà diversa.

	Q1	Q2	Q3	Q4	Q5	%
Studente 1	1	1	1	1	1	100
Studente 2	0	1	1	1	1	80
Studente 3	0	0	1	1	1	60
Studente 4	0	0	0	1	1	40
Studente 5	0	0	0	0	1	20
Studente 6	1	1	0	0	0	40
Df	0,66	0,50	0,50	0,33	0,17	

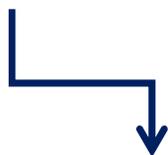
Aspetti critici del livello di difficoltà nella TCT

	Q1	Q2	Q3	Q4	Q5	Q6	%
Studente 1	1	1	1	1	1	0	83
Studente 2	0	1	1	1	1	0	67
Studente 3	0	0	1	1	1	0	50
Studente 4	0	0	0	1	1	0	33
Studente 5	0	0	0	0	1	1	33
Df	0,80	0,60	0,40	0,20	0	0,80	

Indice di discriminatività

Nella composizione di un test o di una prova oggettiva è importante analizzare, tra le altre cose, la capacità che hanno i singoli item di differenziare i soggetti “preparati” da quelli meno preparati.

DOMANDA



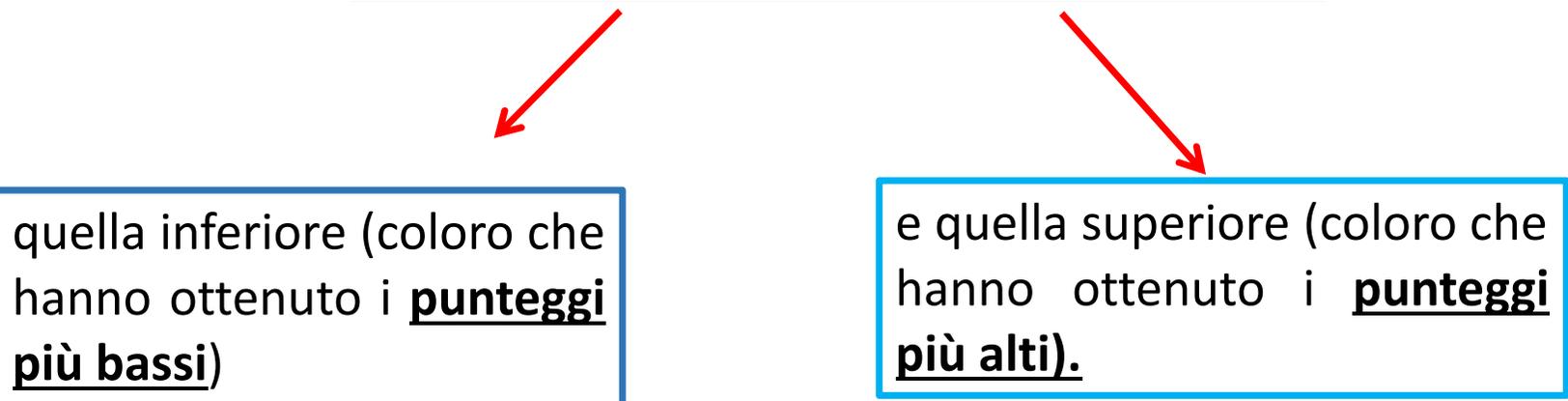
In che misura una domanda posta in un questionario può separare i soggetti che possiedono le abilità che si vogliono verificare con una prova da coloro che non le possiedono?

Uno degli strumenti più utilizzati in tal senso è **l'indice di discriminatività.**

Indice di discriminatività

Il calcolo di questo indice assume che, le risposte giuste date ad una domanda dal gruppo di soggetti più bravi dovrebbero essere più numerose delle risposte esatte date allo stesso quesito dal gruppo che ha conseguito i punteggi più bassi.

È dunque necessario paragonare le due fasce estreme dei risultati



quella inferiore (coloro che hanno ottenuto i punteggi più bassi)

e quella superiore (coloro che hanno ottenuto i punteggi più alti).

Indice di discriminatività

L'**indice di discriminatività/selettività** si calcola nel modo seguente

$$\text{Indice di discriminatività} = (E_s - E_i)/N$$

dove:

E_s = numero delle risposte esatte registrate nell'estremo superiore (numero di allievi della fascia alta che hanno risposto correttamente all'item)

E_i = numero delle risposte esatte registrate nell'estremo inferiore (numero degli allievi della fascia bassa che hanno risposto correttamente all'item)

n = numero dei soggetti che costituiscono ciascun gruppo estremo

L'indice varia da -1 a +1.

Il valore zero indica che l'item non è discriminativo, quando l'estremo superiore risponde meglio il segno è positivo.

Un buon indice di selettività è compreso tra 0,30 e 0,60:
sotto 0,3 → item non è discriminante
sopra 0,6 → item troppo selettivo

Indice di discriminatività

indica quanto l'item riesce a discriminare tra rispondenti con alte abilità (gruppo 1) rispetto a quelli con basse abilità (gruppo 2).

Nello specifico rappresenta la correlazione tra la probabilità di scegliere una data opzione e l'abilità complessiva del rispondente.

In un test a scelta multipla,

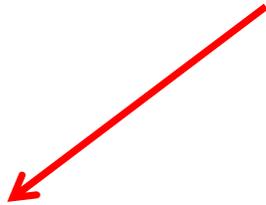
tale legame deve essere negativo per le opzioni di risposta non corrette e positivo solo per quella esatta.

Una domanda è ben formulata se, in media, coloro che rispondono correttamente a quella domanda lo fanno anche a buona parte delle altre che hanno lo stesso livello di difficoltà.

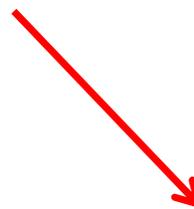
Si considerano funzionanti gli item con un valore $> 0,35$

Indice di discriminatività

L'indice di discriminatività assume valori compresi tra -1 e +1. L'indice discrimina correttamente se assume un valore compreso tra 0,30 e 0,60.



Se per un item l'indice assume valore +1, significa che al quesito hanno risposto correttamente solo studenti di un livello di apprendimento elevato.



Se per un item l'indice assume valori negativi (discriminatività negativa), occorre riflettere sulla formulazione dell'item, in quanto se l'item è troppo fuorviante puo' indurre a risposte casuali.

Esempio

Ad esempio si supponga che il test sia stato svolto da 25 allievi; si individueranno 3 fasce, la prima e la terza di 8 allievi, la seconda di 9 allievi: la prima degli alunni migliori, la terza dei peggiori, la seconda dei mediani.

Si supponga che all'*item* n. 1 abbiano risposto correttamente 6 allievi della fascia alta e 2 allievi della fascia bassa

$$S_1 = (6-2)/8 = 0,50$$

che all'*item* n. 2, rispettivamente 7 e 5 allievi

$$S_2 = (7-5)/8 = 0,25$$

Nel caso da noi ipotizzato l'*item* n. 2 non è discriminante, il n.1 sì.

La discriminatività dei distrattori

Di ogni opzione dovrebbe essere calcolata anche la **discriminatività**, che nel caso **dei distrattori dovrebbe essere negativa**.

Dovrebbe cioè accadere che coloro che hanno punteggio basso nel complesso della prova hanno più facilità di scegliere le risposte errate rispetto a coloro che hanno punteggio alto.

Quando un distrattore ha una discriminatività positiva abbiamo una prova che il tranello insito nel distrattore ha danneggiato i più performanti, e che quindi ha reso meno affidabile la misura operata dal quesito.

Quali quesiti rivedere e perché?

	dom. 1	dom. 2	dom. 3	dom. 4	dom. 5	dom. 6	dom. 7
A	13	10	53	7	9	2	4
B	13	13	5	58	27	4	28
C	45	6	17	3	38	62	9
D	6	48	2	0	2	9	34
Missing	0	0	0	9	1	0	2

	dom. 1	dom. 2	dom. 3	dom. 4	dom. 5	dom. 6	dom. 7
corrette	45	48	53	7	38	62	28
errate	32	29	24	70	39	15	49
Difficoltà	42%	38%	31%	91%	51%	19%	64%
Discrimin.	0,5	0,4	0,4	-0,1	0,7	0,7	0,4

Quali quesiti rivedere e perché?

	Dom. 1	Dom. 2	Dom. 3
A*	24	6	16
B	2	32	11
C	14	5	13
D	10	7	10
diff	0,52	0,88	0,68
discr	0,43	-0,21	0,87

Esempio: prove classi prime matematica

Partecipazione al campionamento: 3 classi prime

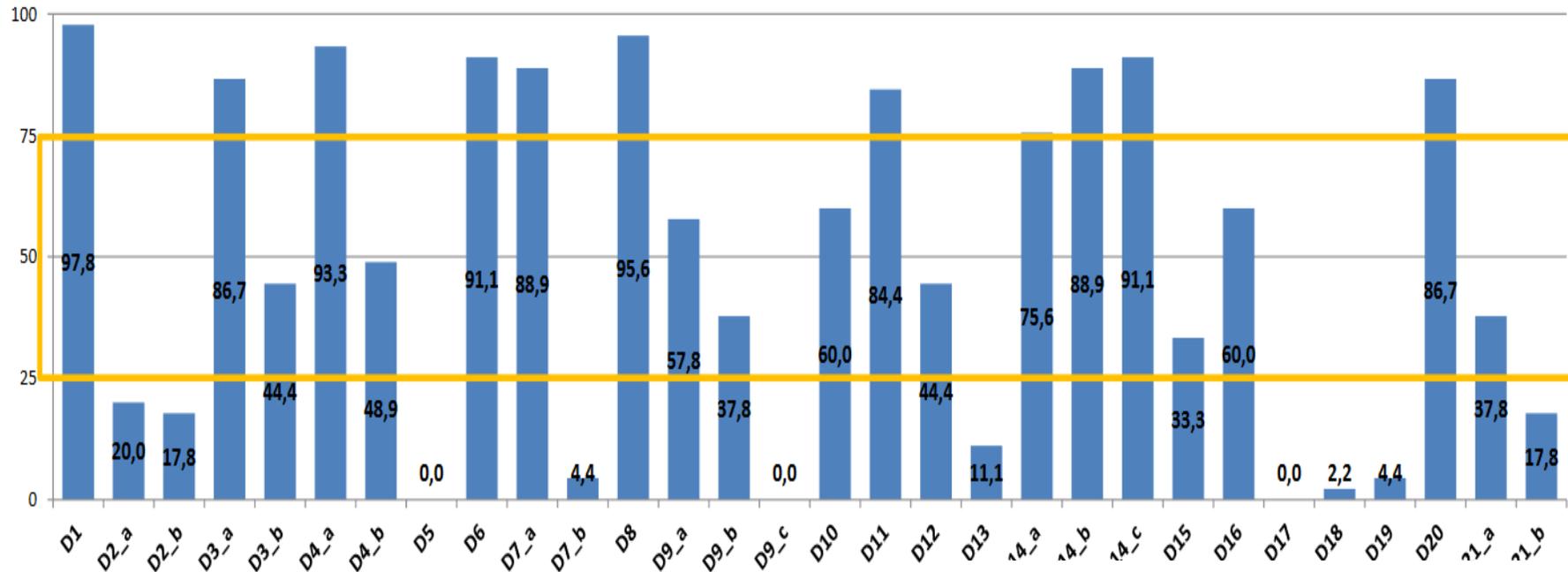
- TOTALE = 45 studenti delle classi prime di 3 classi di 3 scuole diverse.

30 quesiti:

- 15 quesiti a risposta multipla (4 opzioni di risposta A, B, C, D)
- 14 quesiti a risposta libera “breve” → ricodifica omogenea 1/0 dove 1 se corretto e 0 se errato
- 1 quesito complesso D5, composto da 4 item

Statistiche Descrittive		
N studenti partecipanti = 45	# esatte	% esatte
MEDIA	14,8	49,4
MEDIANA	15	50,0
MIN	6	20,0
MAX	18	60,0
RANGE (Max - Min)	12	40,0

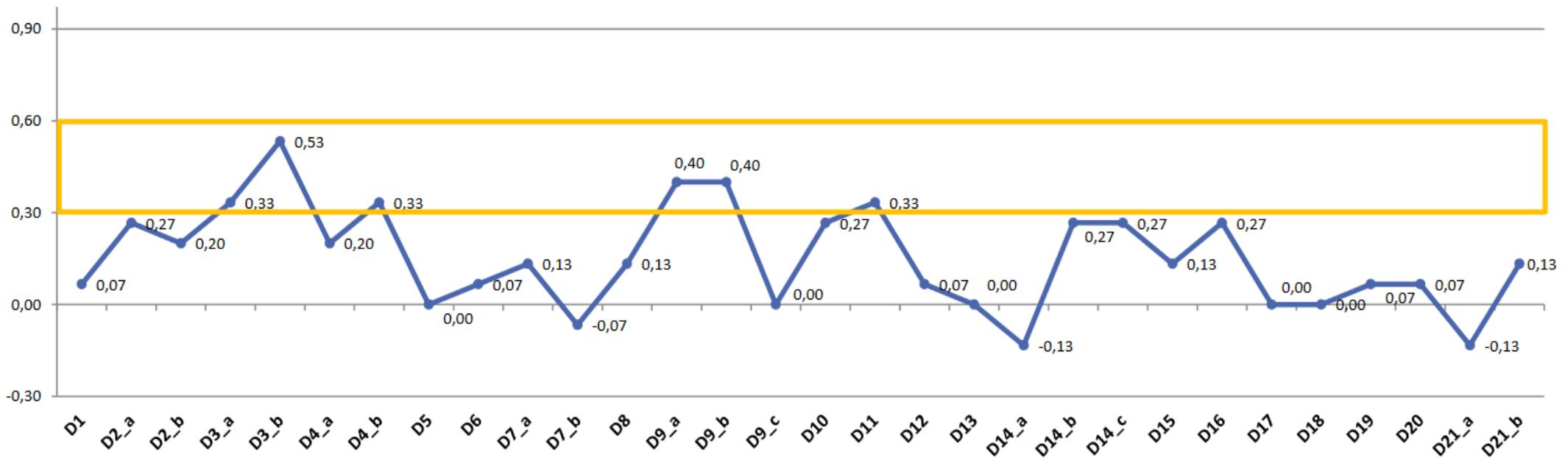
Indice di difficoltà/facilità Studenti classi prime –matematica (N=45)



L'indice di difficoltà /facilità è accettabile se la % di risposte corrette è compresa tra 25% e 75% (oltre 75% l'item è molto facile (FF), sotto il 25% è molto difficile (DD)).

- 9 quesiti con indice di difficoltà/facilità accettabile (6 D, 4 F)
- 10 quesiti troppo facili (FF)
- 10 quesiti troppo difficili (DD)

Indice di discriminatività/selettività Studenti classi prime -matematica (N=45)



L'indice di selettività deve essere compreso tra 0,30 e 0,60
sotto 0,3 → item non è discriminante
sopra 0,6 → item molto selettivo

- 6 quesiti con indice di selettività buono
- Gli altri non discriminanti (di questi 5 sono vicini a 0,30)

Il coefficiente di correlazione punto-biserial

Il coefficiente di correlazione punto-biserial (r_{pb}) è un coefficiente di correlazione utilizzato quando una variabile è dicotomica. La variabile può essere effettivamente dicotomica, come il genere, o può essere dicotomizzata artificialmente. Il coefficiente di correlazione punto-biserial è che la correlazione di Pearson tra l'item e il punteggio totale al test. La sua espressione è data nella seguente

$$r_{pb} = \frac{M_1 - M_0}{\sigma_n} \sqrt{\frac{n_1 n_0}{n^2}}$$



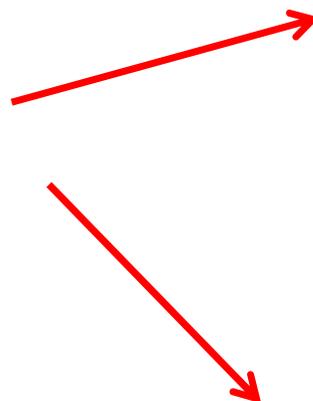
rappresenta la deviazione standard dell'intera popolazione

M1 è il **valore medio della variabile continua X** per tutti i dati del **gruppo 1**, rappresenta quindi la media dei punteggi dei soggetti che hanno risposto in maniera corretta all'item;

M0 è il **valore medio delle X per tutti i dati del gruppo 2**; n_1 è la numerosità del gruppo 1; n_0 è la numerosità del gruppo 2 ed n è la dimensione totale del campione.

Il coefficiente di correlazione punto-biserial

Il coefficiente di correlazione punto-biserial stabilisce una correlazione tra due variabili



una dicotomica che consiste nell'aver risposto in maniera esatta o errata ad uno specifico item

l'altra continua che consiste nel punteggio complessivo ottenuto da coloro che hanno risposto esattamente all'item

Il coefficiente di correlazione punto-biseriale

Indice di discriminazione

Item 4

item:4 (Flps3_spedizione_bici)
Cases for this item 105 Discrimination 0.41
Item Threshold(s): 0.22 Weighted MNSQ 1.21 **Fit**
Item Delta(s): 0.22

Label	Score	Count	% of tot	Pt Bis	t (p)	PV1Avg:1	PV1 SD:1
0	0.00	17	16.19	-0.32	-3.45 (.001)	-0.82	1.05
1	1.00	50	47.62	0.41	4.53 (.000)	0.47	1.11
2	0.00	38	36.19	-0.18	-1.83 (.070)	-0.15	1.30

Risposta corretta (pointing to Score 1.00)

Percentuale di risposte corrette o indice di facilità (pointing to % of tot 47.62)

Correlazione punto biseriale (pointing to Pt Bis 0.41)

Indice di discriminazione (pointing to Discrimination 0.41)

Fit (pointing to Weighted MNSQ 1.21)

Il coefficiente di correlazione punto-biserial

Il coefficiente di correlazione punto-biserial stabilisce una correlazione tra tutte le risposte date ad un quesito e tutti i punteggi grezzi degli allievi

Differisce dalla discriminatività che è un semplice confronto tra gli estremi

Varia da -1 a +1 e si interpreta

- ✓ <0 = DA SCARTARE
- ✓ <0.19 = RIVEDERE IL QUESITO
- ✓ $0.20-0.29$ = QUESITI ACCETTABILI
- ✓ $0.30-0.39$ = QUESITI BUONI
- ✓ >0.40 = QUESITI OTTIMI

Il coefficiente di discriminatività e la correlazione punto-biseriale

	Quesito 1	Quesito 2	Quesito 3
Estremo superiore	5	4	6
Estremo inferiore	1	5	0
Numero studenti per fascia	6	6	6
Discriminatività	0,7	-0,2	1
Correlazione punto biseriale	0,6	-0,2	0,7

La correlazione punto-biserial in Excel



Esempio

Students	Items									
	1	2	3	4	5	6	7	8	9	10
Kid-A	1	1	1	1	1	1	1	1	0	1
Kid-B	1	1	1	1	1	1	1	0	1	0
Kid-C	1	1	1	1	1	1	0	1	0	0
Kid-D	1	1	1	1	1	0	1	0	1	0
Kid-E	1	1	1	1	1	1	0	1	0	0
Kid-F	1	1	1	0	1	0	0	0	0	0
Kid-G	1	1	0	1	0	1	0	0	0	0
Kid-H	1	0	1	0	1	0	0	0	0	0
Kid-I	0	1	1	0	0	0	0	0	0	0

	A	B	C	D	E	F	G	H	I	J	K	L
1	Items	1	2	3	4	5	6	7	8	9	10	Student Total Score
2	Students											
2	Kid-A	1	1	1	1	1	1	1	1	0	1	9
3	Kid-B	1	1	1	1	1	1	1	0	1	0	8
4	Kid-C	1	1	1	1	1	1	0	1	0	0	7
5	Kid-D	1	1	1	1	1	0	1	0	1	0	7
6	Kid-E	1	1	1	1	1	1	0	1	0	0	7
7	Kid-F	1	1	1	0	1	0	0	0	0	0	4
8	Kid-G	1	1	0	1	0	1	0	0	0	0	4
9	Kid-H	1	0	1	0	1	0	0	0	0	0	3
10	Kid-I	0	1	1	0	0	0	0	0	0	0	2
11	Item total	8	8	8	6	7	5	3	3	2	1	

Note: Callouts in the original image indicate formulas: '=sum(B2:K2)' for the Student Total Score and '=sum(B2:B10)' for the Item total.

La correlazione punto-biserial in Excel

Esempio



`=L2-B2`

		M	N	O	P	Q	R	S	T	U	V
1	Students	total-item1	total-item2	total-item3	total-item4	total-item5	total-item6	total-item7	total-item8	total-item9	total-item10
2	Kid-A	8	8	8	8	8	8	8	8	9	8
3	Kid-B	7	7	7	7	7	7	7	8	7	8
4	Kid-C	6	6	6	6	6	6	7	6	7	7
5	Kid-D	6	6	6	6	6	7	6	7	6	7
6	Kid-E	6	6	6	6	6	6	7	6	7	7
7	Kid-F	3	3	3	4	3	4	4	4	4	4

`=CORREL(B2:B10, M2:M10)`

		M	N	O	P	Q	R	S	T	U	V
		total-item1	total-item2	total-item3	total-item4	total-item5	total-item6	total-item7	total-item8	total-item9	total-item10
12	Point-Biserial	0.46	0.29	0.12	0.73	0.49	0.49	0.59	0.46	0.26	0.40

Analisi dell'affidabilità di un test

Lo studio dell'affidabilità o attendibilità di un questionario riguarda l'analisi della coerenza tra i punteggi ottenuti dalla somministrazione di una prova se questa viene somministrata in momenti successivi.

ATTENDIBILITA'



è definibile come il rapporto fra la componente sistematica e la variabilità totale della misura

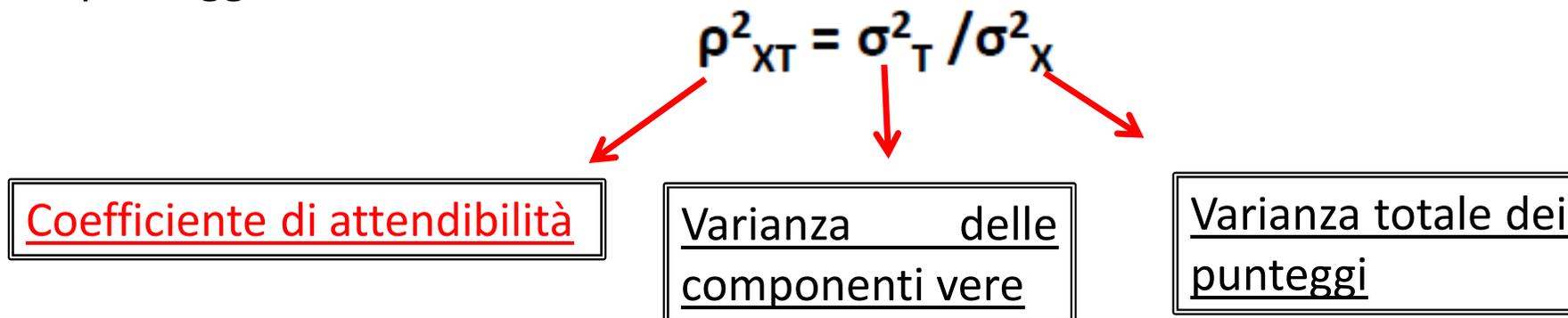
Affinché una misura sia valida, si richiede che essa misuri sistematicamente qualcosa. L'attendibilità rappresenta la precisione di una misura (ciò che nella misura non è errore).

Analisi dell'affidabilità di un test

La **variazione dei punteggi** può essere dovuta a componenti vere e a componenti di errore. Si può allora affermare che la varianza totale è pari alla somma della varianza delle componenti vere e di quella delle componenti di errore

$$\sigma^2_X = \sigma^2_T + \sigma^2_E$$

L'attendibilità di un test può essere definita come la proporzione della varianza vera rispetto alla varianza totale di una distribuzione dei punteggi



Analisi dell'affidabilità di un test

L'attendibilità di un test può essere scritta come

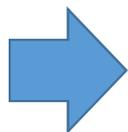
$$\rho_{XT}^2 = 1 - \frac{\sigma_E^2}{\sigma_X^2}$$

da cui si ricava che

$$\sigma_E = \sqrt{\sigma_X^2(1 - \rho_{XT}^2)}$$

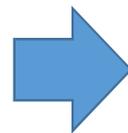
Si comprende la relazione inversa che lega l'affidabilità all'errore

$$\text{se } \rho_{XT}^2 = 1$$



Tutta la variazione dei valori osservati è attribuibile ai punteggi veri

$$\text{se } \rho_{XT}^2 = 0$$



Tutta la variazione dei valori osservati è attribuibile all'errore

Interpretazione dell'affidabilità

Varia tra:

- 0, il punteggio è composto solo da errore
- 1, il punteggio osservato è vero

da 0 a 1 si esprimono valori intermedi di attendibilità

L'attendibilità è buona se superiore a 0,70 – 0,80; è massima se è 1.

Analisi dell'affidabilità di un test

Poiché la varianza dei punteggi osservati può essere considerata in termini di somma della varianza del punteggio vero e della varianza dell'errore, possiamo esprimere

$$\rho_{XT}^2 = \frac{\sigma_T^2}{\sigma_X^2}$$

$$\rho_{XX'} = \frac{\sigma_T^2}{\sigma_X^2} = \frac{\sigma_T^2}{\sigma_T^2 + \sigma_E^2}$$

da cui si deduce che il quadrato della correlazione tra i punteggi osservati e i punteggi veri è uguale alla correlazione tra i punteggi osservati di due misurazioni parallele.

Quindi l'affidabilità dei punteggi diventa maggiore se la percentuale di varianza dell'errore è piccola e viceversa.

La radice quadrata dell'indice di affidabilità rappresenterà la correlazione tra i punteggi veri e i punteggi osservati.

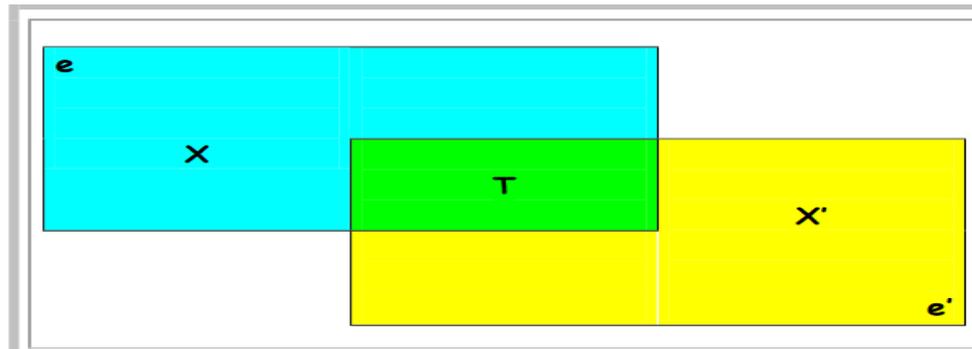
Formalmente si ha....

La correlazione tra le due misure ripetute x e x' consente di stimare l'affidabilità di una misura.

Le due misure ripetute (x e x') sono dette parallele se presentano:

1. uguali punteggi veri attesi per ciascun oggetto,
2. uguali varianze dell'errore o, in modo equivalente, errori standard di misurazione, ovvero se risultano vere le seguenti relazioni

$$X = T + E \quad X' = T' + E' \quad T = T' \quad \sigma_E^2 = \sigma_{E'}^2$$



Formalmente si ha....

La correlazione tra misure parallele può essere espressa in funzione dell'errore, del punteggio vero e del punteggio osservato:

$$\rho_{XX'} = \frac{\sigma_{XX'}}{\sigma_X \sigma_{X'}} = \frac{\sigma_{(T+E)} \sigma_{(T+E')}}{\sigma_X \sigma_{X'}} = \frac{\sigma_T^2 + \sigma_{TE} + \sigma_{TE'} + \sigma_{EE'}}{\sigma_X \sigma_{X'}}$$

ASSUMIAMO CHE

1. GLI ERRORI NON SONO CORRELATI NE' TRA LORO, NE' CON I PUNTEGGI VERI
2. LE DEVIAZIONI STANDARD DELLE MISURE PARALLELE SONO UGUALI

La correlazione tra misure parallele è uguale al rapporto tra varianza dei punteggi veri e varianza dei punteggi osservati:

$$\rho_{XX'} = \frac{\sigma_T^2}{\sigma_X^2}$$

Formalmente si ha....

$$\rho_{XX'} = \frac{\sigma_T^2}{\sigma_X^2}$$

Tale risultato è importante in quanto consente di esprimere la varianza del punteggio vero (inosservabile) in termini di $\rho_{XX'}$ e σ_X^2 , entrambi osservabili:

$$\sigma_T^2 = \sigma_X^2 \rho_{XX'}$$

ovvero la varianza del punteggio vero è uguale al prodotto tra la varianza osservata e la correlazione tra misure parallele. Ricordando che l'affidabilità è stata definita come $\rho_{XT}^2 = \frac{\sigma_T^2}{\sigma_X^2}$, ne segue che la stima dell'affidabilità non è altro che la correlazione tra misure parallele:

$\rho_{XT}^2 = \frac{\sigma_T^2}{\sigma_X^2} = \frac{\sigma_X^2 \rho_{XX'}}{\sigma_X^2} = \rho_{XX'}$. Tale risultato rappresenta un importante risultato nel tentativo di stimare l'affidabilità delle misure empiriche.

L'attendibilità

Il coefficiente ρ_{XT}^2



È il coefficiente di attendibilità del test

RELIABILITY

$\rho_{XX'}$ → È l'indice di affidabilità
Serve a valutarla empiricamente

L'indice di affidabilità tuttavia fornisce indicazioni circa il livello di attendibilità dell'intero questionario, ma non dà alcuna informazione in merito all'affidabilità dei singoli item; questo limite viene superato, come vedremo, dall'IRT.

Analisi dell'affidabilità di un test

La Teoria Classica dei Test studia l'attendibilità di un questionario in modo non univoco.

Il coefficiente di attendibilità è stimato in vari modi:

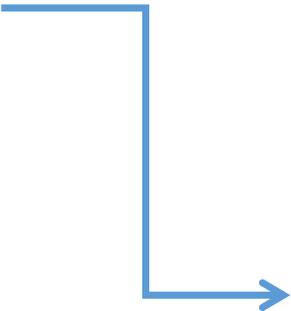
- 1) come correlazione tra i punteggi conseguiti da un gruppo di soggetti in due metà dello stesso test (metodo dello Split-Half);
- 2) come correlazione tra le due distribuzioni di punteggi ottenute applicando due volte uno stesso test ad uno stesso gruppo di soggetti (metodo test-retest);
- 3) come correlazione tra le due serie di punteggi ottenute somministrando allo stesso gruppo di soggetti due forme parallele dello stesso test (parallel form);
- 4) come studio della coerenza o omogeneità tra gli item (alpha di Cronbach).

Lo split-half

Si somministra il test in un unico tempo T1.

Si divide il test a metà e si considerano le due metà come forme parallele (stessa media e stessa dev. St.)

La suddivisione delle variabili in due gruppi deve essere effettuato in modo che essi risultino omogenee tra loro.



I questionari per la valutazione degli apprendimenti sono costruiti in modo tale che i diversi item abbiano un livello di difficoltà crescente. Di conseguenza una suddivisione a metà in relazione all'ordinamento potrebbe determinare che nella parte più semplice si ottengano risultati migliori rispetto alla parte più complessa. Per ovviare a questa problematica si adotta il metodo "odd-even".

Lo split-half

Prendiamo i seguenti dati di un test di 10 item somministrato a 11 soggetti:

Item	1	2	3	4	5	6	7	8	9	10
Sogg.										
1	1	0	1	1	0	1	0	1	1	1
2	0	0	1	0	1	1	0	1	0	0
3	1	1	0	1	1	1	1	0	1	1
4	0	0	1	0	0	1	0	1	0	0
5	1	0	1	1	0	1	0	0	1	1
6	1	1	0	0	1	0	1	1	0	0
7	0	0	1	0	0	0	1	0	0	0
8	1	1	1	1	0	1	0	0	0	0
9	0	1	0	0	1	1	1	0	0	1
10	1	0	0	1	0	0	1	0	0	1
11	1	1	1	0	1	0	0	1	1	0

Sommiamo gli item pari e dispari

Sogg.	Item Dispari	Item Pari
1	3	4
2	2	2
3	4	4
4	1	2
5	3	3
6	3	2
7	2	0
8	2	3
9	2	3
10	2	2
11	4	2
Media	2,55	2,45
Dev.St.	0,89	1,08

La correlazione tra le due metà=0.40

Interpretazione coefficiente split-half

L'attendibilità dipende molto dalla lunghezza del test.



la correlazione split-half è una sottostima dell'attendibilità.
Infatti la divisione del test a metà ne dimezza la lunghezza.

Ci sono delle formule che permettono di correggere tale sottostima.

Formula di Spearman-Brown

$$\rho_{XX'}^* = \frac{n\rho_{xx'}}{1 + (n-1)\rho_{xx'}}$$

n= numero delle volte che il test viene “allungato” o “abbreviato”.
Si ottiene come rapporto tra numero degli item iniziali su numero item finali.

Attendibilità dell'intero test

Attendibilità iniziale

Applicazione al problema precedente

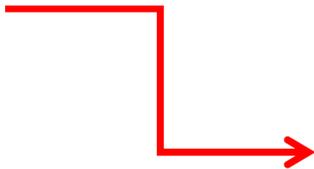
$$\rho_{xx'} = 0.40 \quad n = \frac{10}{5} = 2$$

$$\rho_{XX'}^* = \frac{2 \times 0.40}{1 + (2 - 1) \times 0.40} = \frac{0.80}{1.40} = 0.58$$

Il metodo del Test-Retest

Si somministra il **test al tempo T1 ed al tempo T2** e si calcola **la correlazione** tra i punteggi.

Questo metodo non necessita di ulteriori specificazioni. Basta saper calcolare la **r di Pearson** tra **due serie di punteggi**.



Il coefficiente di attendibilità ottenuto esprime, quindi, il grado di stabilità del test e di generalizzabilità dei risultati in caso di somministrazioni diverse.

Quanto più alto risulterà il coefficiente di attendibilità, tanto minore sarà l'influenza delle variabili accidentali sui punteggi.

Il metodo del Test-Retest

	PG T1	PG T2
ss1	11	12
ss2	15	14
ss3	17	14
ss4	20	19
ss5	20	21
ss6	25	27
ss7	22	18
ss8	21	24
ss9	34	31
ss10	38	36
ss11	40	37
Media	23,91	23,00
Dev St.	9,03	8,39

$$Cov (PG T1, PG T2) = 73.45$$

$$Dev.St.(PG T1) = 9.03$$

$$Dev.St.(PG T2) = 8.39$$

$$\rho_{xx'} = \frac{Cov (PG T1, PG T2)}{Dev.St.(PG T1) \times Dev.St.(PG T2)}$$

$$73.45 / 9.03 \times 8.39 =$$

$$73.45 / 75.74 = .96$$

Interpretazione coefficiente test-retest

- Buoni coefficienti test-retest dovrebbero superare .80 (livello piuttosto esigente).
- Il coefficiente test-retest si riduce all'aumentare del tempo trascorso fra le rilevazioni.
- Il coefficiente test-retest è interpretabile se si assume che il concetto misurato non si modifichi nel tempo

Parallel form

Si somministrano due versioni equivalenti del test (stessa media e stessa dev. St.) Quindi si calcola la correlazione tra i le due forme come stima dell'attendibilità test-retest.

OSSERVAZIONE:

- Se le due versioni del test sono state somministrate a distanza di tempo, tale coefficiente di correlazione diventa un indice sia della stabilità del tempo, sia della coerenza delle risposte date a diversi campioni di prove che presentano contenuti simili.
- Se invece le due forme vengono somministrate consecutivamente, il coefficiente di correlazione trovato esprimerà il grado di attendibilità tra le due versioni, ma non tra i due tempi di applicazione.

Parallel form

Il problema principale è quello di verificare che le due forme siano effettivamente parallele. Ciò significa verificare che le due forme abbiano la stessa media e la stessa varianza.

Prendiamo i seguenti dati:

	T1 Forma A	T2 Forma B
ss1	11	12
ss2	15	14
ss3	17	14
ss4	20	19
ss5	20	21
ss6	25	27
ss7	22	18
ss8	21	24
ss9	34	31
ss10	38	36
ss11	40	37
Media	23,91	23,00
Dev St.	9,03	8,39

t-test sulle due medie

$$t_{sp}=1.3; Gdl=10; t_{cr}=2,23$$

Le medie sono uguali.

Fare Test F sulle due varianze

$$F_{sp}=1.15, gdl = 10,10$$
$$F_{cr}=2.97$$

Le varianze sono uguali.

Correlazione



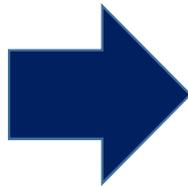
$$\rho_{xx'} = 0.96$$

L'alpha di Cronbach

L'attendibilità di un test può essere studiata in funzione della **coerenza od omogeneità degli item all'interno di un test.**

L'attendibilità viene quindi stimata utilizzando l'informazione contenuta negli item. Il coefficiente più utilizzato è l'alpha di Cronbach.

$$\alpha = \frac{k}{k-1} \left(1 - \frac{\sum_{i=1}^k \sigma_i^2}{\sigma_X^2} \right)$$



Si utilizza per item dicotomici
o politomici

in cui:

k = numero di item

δ_i^2 = varianza di ciascun item

δ_x^2 = varianza totale del test

Esprime una misura del peso relativo della variabilità associata agli item rispetto alla variabilità associata alla loro somma.

alpha di Cronbach: esempio dati dicotomici

Item Sogg.	1	2	3	4	5	6	7	8	9	10	Totale
1	1	0	1	1	0	1	0	1	1	1	7
2	0	0	1	0	1	1	0	1	0	0	4
3	1	1	0	1	1	1	1	0	1	1	8
4	0	0	1	0	0	1	0	1	0	0	3
5	1	0	1	1	0	1	0	0	1	1	6
6	1	1	0	0	1	0	1	1	0	0	5
7	0	0	1	0	0	0	1	0	0	0	2
8	1	1	1	1	0	1	0	0	0	0	5
9	0	1	0	0	1	1	1	0	0	1	5
10	1	0	0	1	0	0	1	0	0	1	4
11	1	1	1	0	1	0	0	1	1	0	6
Media	0,64	0,45	0,64	0,45	0,45	0,64	0,45	0,45	0,36	0,45	5,00
Dev.St.	0,48	0,50	0,48	0,50	0,50	0,48	0,50	0,50	0,48	0,50	1,65
Varianza	0,23	0,25	0,23	0,25	0,25	0,23	0,25	0,25	0,23	0,25	2,73

Il coefficiente α di Cronbach

$$\alpha = \frac{K}{K-1} \times \left(1 - \frac{\sum \sigma_i^2}{\sigma_{totale}^2} \right) =$$

$$= \frac{10}{9} \times \left(1 - \frac{.23 \times 4 + .25 \times 6}{2.73} \right) =$$

$$= \frac{10}{9} \times \left(1 - \frac{2.42}{2.73} \right) = .1261$$

Alpha di Cronbach: interpretazione

l' Alpha di Cronbach:

dipende dalla media delle intercorrelazioni tra tutti gli item del test, e dalla relazione di ogni item del test con il punteggio totale.

Nella prassi si valuta così:

α di Cronbach	Coerenza interna del test
$\alpha \geq 0.9$	Eccellente
$0.8 \leq \alpha < 0.9$	Buono
$0.7 \leq \alpha < 0.8$	Accettabile
$0.6 \leq \alpha < 0.7$	Discutibile
$0.5 \leq \alpha < 0.6$	Povero
$0. \alpha < 0.5$	Inaccettabile

Considerazioni pratiche

Per meglio valutare l'Alpha di Cronbach, bisogna ricordare che essa **dipende da due fattori:**

1. **intercorrelazioni tra gli item**
2. **lunghezza della scala (numero di item)**

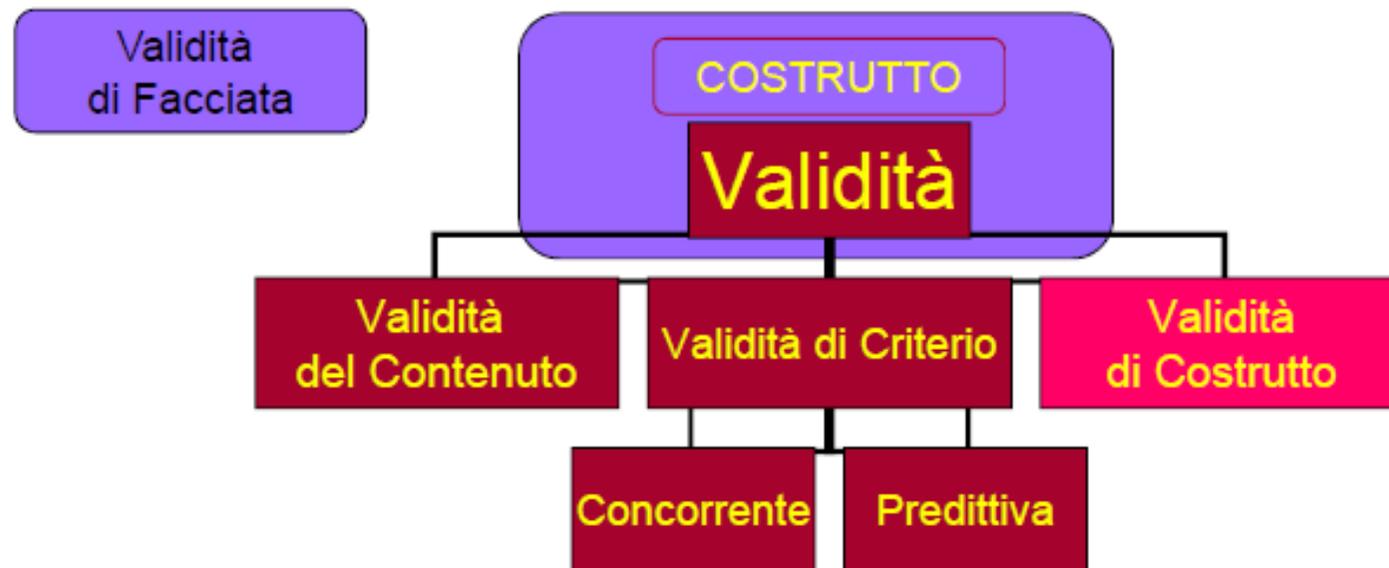
Infatti, a parità di condizioni, all'aumentare del numero degli item, aumenta il valore del coefficiente di attendibilità.

Perciò, per esempio, basterebbe aumentare il numero di item di una scala che ha un coefficiente sufficiente per ottenerne uno buono.

Analisi della validità di un test

- Una misura è valida quando misura ciò che intende misurare.
- Si tratta di *“un giudizio complessivo della misura in cui prove empiriche e principi teorici supportano l’adeguatezza e l’appropriatezza delle conclusioni basate su punteggi al test”* (Messick, 1989)
- La validità si articola in diverse sfumature, e in diverse modalità di giudicare, raggiunte o meno, differenti forme di validità

Tassonomia tradizionale



I diversi aspetti della validità

Validità interna:

gli item sono misure del costrutto sufficientemente correlate tra loro. L'analisi fattoriale consente di verificare tale aspetto.

Validità di criterio

grado di corrispondenza tra una misura ed un criterio di riferimento. Si distingue in:



Validità concorrente

misura e criterio sono rilevati nello stesso momento

Validità predittiva

il criterio è misurato successivamente alla misura.

Validità di criterio

Esempio di validità concorrente:

Correlazione tra quoziente intellettivo misurato tramite un test di intelligenza e la risoluzione di un determinato problema di una certa complessità cognitiva (criterio)

Esempio di validità predittiva:

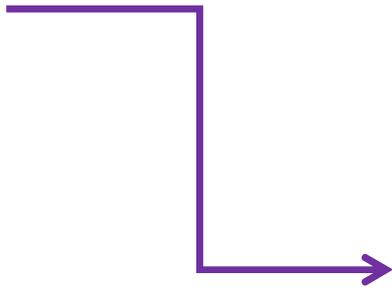
Correlazione tra il punteggio ottenuto ad un test attitudinale nella fase di selezione del personale per una certa azienda ed il successo lavorativo (ad es., velocità della carriera interna) ottenuto negli anni successivi (criterio).

Analisi Fattoriale

la dimensionalita' di un test

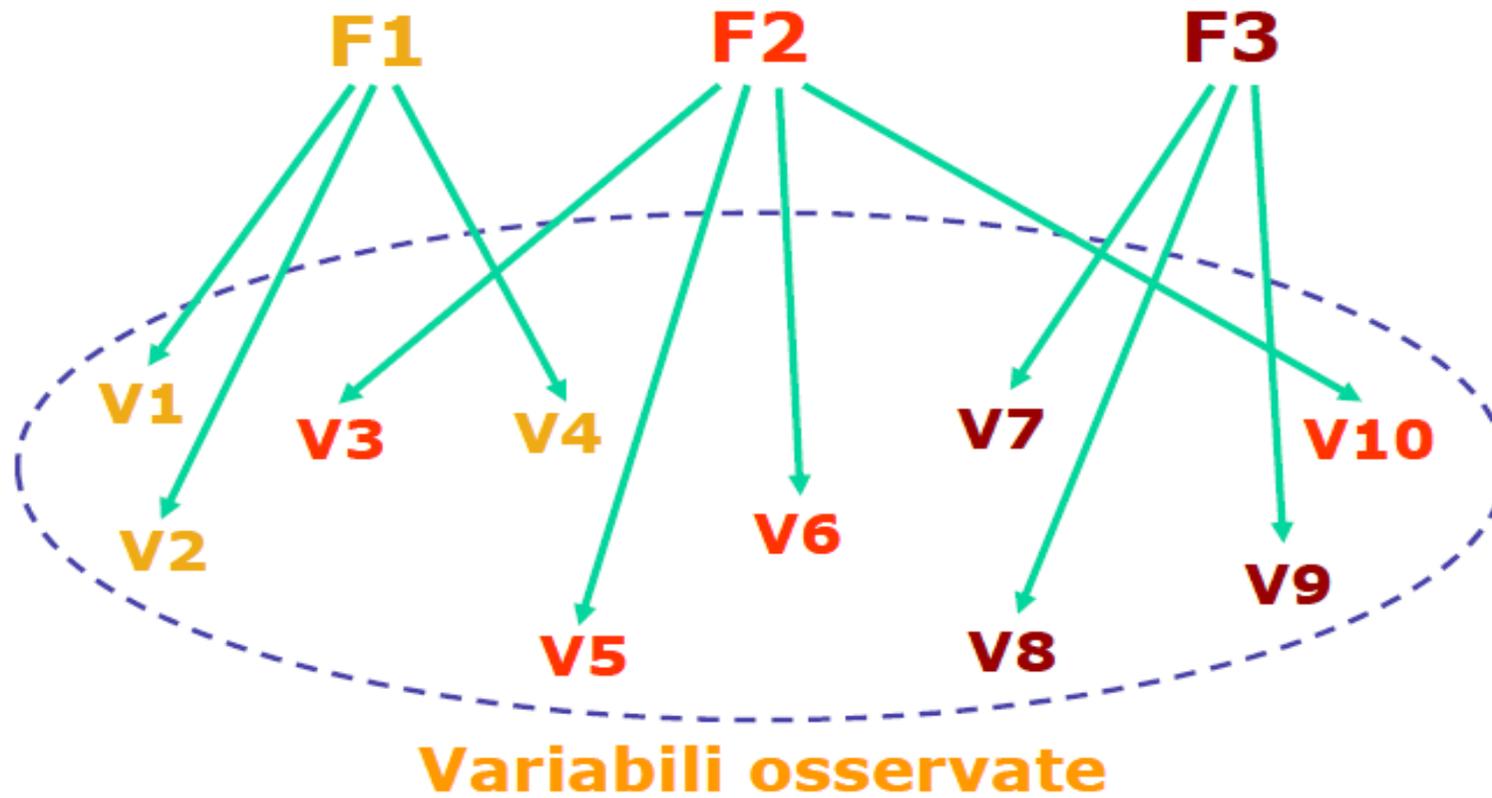
Studiare la dimensionalità di un test significa studiare le dimensioni latenti ad un insieme di item e ciò viene effettuato tramite **l'Analisi Fattoriale**.

La dimensione latente può essere considerata un costrutto teorico ipotetico, idealmente corrispondente al costrutto teorico inizialmente ipotizzato.



L'analisi fattoriale viene utilizzata per identificare un numero di fattori latenti (dimensioni, tratti, componenti) che spieghino le correlazioni tra le variabili osservate (indicatori, item) **in modo parsimonioso**.
I fattori perciò sono sempre in numero molto minore rispetto agli item analizzati.

Fattori latenti



Correzione per guessing

La correzione per il guessing scompone il numero totale di risposte corrette in due componenti:

1. le risposte corrette dovute alle conoscenze del soggetto
2. e quelle che risultano corrette come effetto del caso.

Individui con lo stesso numero di risposte corrette, ma un diverso numero di errate e/o omesse, non otterranno lo stesso punteggio corretto per "guessing."

$$P_g = C - \frac{E}{(K - 1)}$$

Dove:

- P_g è il PUNTEGGIO CORRETTO PER GUESSING
- C indica il numero delle risposte CORRETTE
- E indica il numero delle risposte ERRATE
- K indica le possibili alternative di risposta

Correzione per guessing

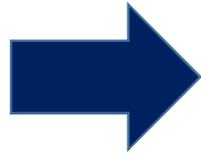
Esempio

Consideriamo il caso di Tonio e Riccardo, che hanno sostenuto la stessa prova d'esame.
Il questionario è composto da 40 item a 5 alternative di risposta.

Tonio risponde correttamente a 32 risposte, ne sbaglia 2 e ne omette 6; Riccardo ottiene 32 corrette, 6 errate e 2 omesse, quindi sommando semplicemente le risposte corrette entrambi gli studenti otterrebbero un punteggio grezzo di 32.

I loro punteggi corretti saranno, invece:

$$X_{\text{Tonio}} = 32 - \frac{2}{(5-1)} = 31,50 \quad X_{\text{Riccardo}} = 32 - \frac{6}{(5-1)} = 30,50$$



Per effetto della correzione per guessing a parità di risposte esatte, Tonio ottiene un punteggio superiore rispetto a Riccardo nella graduatoria finale dei punteggi grezzi.

Correzione per Guessing

-Esempio

Un questionario può essere costituito da item di diversa tipologia risposte a scelta multipla, risposte dicotomiche, risposte di tipo cloze, ecc..

Per fare un semplice esempio concreto, supponiamo che un questionario sia composto da 5 domande Vero/Falso (gruppo A), 4 domande a tre alternative di risposta (gruppo B) e 4 domande di tipo cloze a quattro alternative di completamento (gruppo C). Tonio ottiene 3 risposte corrette al gruppo A di domande, 2 al gruppo B e 3 al gruppo C. Il suo punteggio totale sarebbe quindi pari a 8. Il suo punteggio totale corretto, invece, per ogni blocco di domande sarà:

$$X_{\text{Tonio,A}} = 3 - \frac{2}{(2-1)} \quad X_{\text{Tonio,B}} = 2 - \frac{2}{(3-1)}$$

$$X_{\text{Tonio,C}} = 3 - \frac{1}{(4-1)}$$

da cui un punteggio totale corretto è uguale a:

$$X_{\text{Tot}} = X_{\text{Tonio,A}} + X_{\text{Tonio,B}} + X_{\text{Tonio,C}} = 2 + 1 + 2,66 = 5,66$$

ANALISI DI UN TEST SECONDO LA TCT:

Caso di studio

Un questionario per la valutazione dell'abilità di comprensione del testo è stato somministrato a tutti gli studenti delle classi quarte di una Scuola Primaria.

1. Hanno compilato il questionario 84 studenti
2. Prova costituita da 23 item a scelta multipla

Analisi di facilità degli item

Proportions for each level of response:

	0	1	logit
X1t	0.1548	0.8452	1.6977
X2t	0.0595	0.9405	2.7600
X3t	0.0119	0.9881	4.4188
X4t	0.2738	0.7262	0.9754
X5t	0.0476	0.9524	2.9957
X6t	0.3571	0.6429	0.5878
X7t	0.0595	0.9405	2.7600
X8t	0.1786	0.8214	1.5261
X9t	0.0357	0.9643	3.2958
X10t	0.1548	0.8452	1.6977
X11t	0.2738	0.7262	0.9754
X12t	0.4762	0.5238	0.0953
X13t	0.8690	0.1310	-1.8926
X14t	0.4167	0.5833	0.3365
X15t	0.3095	0.6905	0.8023
X16t	0.3929	0.6071	0.4353
X1i	0.4881	0.5119	0.0476
X2i	0.2024	0.7976	1.3715
X3i	0.2619	0.7381	1.0361
X4i	0.5952	0.4048	-0.3857
X5i	0.2976	0.7024	0.8587
X6i	0.2262	0.7738	1.2299
X7i	0.2857	0.7143	0.9163

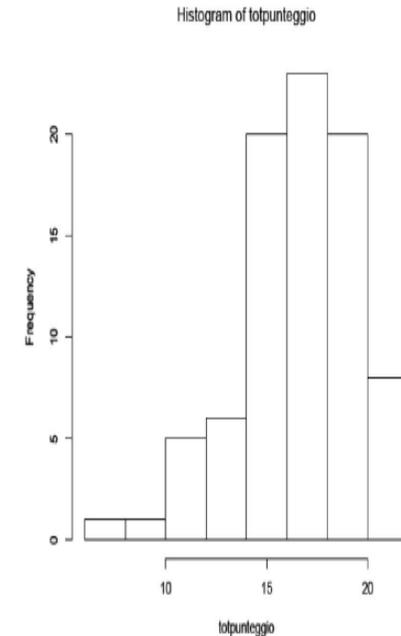
troppo facile

troppo difficile

Analisi della correlazione punto biseriale e dell'alpha di Cronbach

La correlazione punto-biseriale
Varia da -1 a 1

	Point Biserial correlation with Total Score:	Cronbach's alpha all Items:
X1t	0.3883	0.6587
X2t	0.3531	0.6624
X3t	0.1878	0.6716
X4t	0.1241	0.6891
X5t	0.4511	0.6577
X6t	0.4664	0.6513
X7t	0.3531	0.6624
X8t	0.3395	0.6637
X9t	0.4080	0.6613
X10t	0.3277	0.6644
X11t	0.4438	0.6536
X12t	0.2917	0.6748
X13t	0.0294	0.6878
X14t	0.5486	0.6400
X15t	0.4890	0.6481
X16t	0.1561	0.6898
X1l	0.3687	0.6650
X2l	0.5796	0.6379
X3l	0.5285	0.6431
X4l	0.4138	0.6587
X5l	0.3460	0.6657
X6l	0.3919	0.6591
X7l	0.2000	0.6817



Analisi dei distrattori

\$`item_ X4t` score.level	\$`item_ X13t` score.level	\$`item_ X16t` score.level	\$`item_ X71` score.level
resp. lower middle upper	resp. lower middle upper	resp. lower middle upper	resp. lower middle upper
*A 0.703 0.619 0.846	A 0.378 0.238 0.269	A 0.054 0.000 0.000	A 0.135 0.190 0.077
B 0.027 0.143 0.000	*B 0.135 0.190 0.077	B 0.135 0.048 0.000	*B 0.676 0.667 0.808
C 0.162 0.095 0.154	C 0.351 0.524 0.615	C 0.243 0.381 0.308	C 0.135 0.048 0.077
D 0.108 0.143 0.000	D 0.135 0.048 0.038	*D 0.568 0.571 0.692	D 0.054 0.095 0.038

13t. I tempi dei fatti narrati sono:

- A.** Recenti e precisati.
- B.** Passati e non precisati.
- C.** Recenti e non precisati.
- D.** Passati e precisati.

Limiti della TCT

I limiti della TCT riguardano l'impossibilità

- di separare le caratteristiche delle persone da quelle degli item;
- di determinare, nella pratica, indici per la verifica dell'affidabilità dei test;
- di studiare il comportamento di un singolo individuo nei confronti di un singolo item in quanto si limita a fornire statistiche a livello generale dei test.