

Biostatistica

23 novembre 2025



# Indice

<b>I</b>	<b>Misc</b>	<b>9</b>
<b>1</b>	<b>Associazione tra variabili</b>	<b>11</b>
1.1	Entrambe quantitative . . . . .	11
1.2	Una quantitativa e una ordinale . . . . .	11
1.3	Entrambe ordinali . . . . .	11
1.3.1	Kendall $\tau$ (tau) . . . . .	11
1.3.1.1	Prima versione: $\tau_A$ . . . . .	11
1.3.1.2	Seconda versione: $\tau_B$ . . . . .	12
1.3.1.3	Terza versione . . . . .	12
1.3.2	Goodman e Kruskal $\gamma$ . . . . .	13
1.4	Una ordinale e una nominale . . . . .	13
<b>II</b>	<b>Health economics</b>	<b>15</b>
<b>2</b>	<b>Introduzione</b>	<b>17</b>
2.1	Economia sanitaria . . . . .	17
2.2	Sistema sanitario e valutazione economica . . . . .	17
2.3	Processo di valutazione economica in sanità . . . . .	18
2.4	Cost utility shit . . . . .	20
2.5	Comparazione di costi ed effetti . . . . .	21
<b>3</b>	<b>Outcomes</b>	<b>23</b>
3.1	Costi . . . . .	23
3.1.1	Attualizzazione . . . . .	23
3.2	Efficacia . . . . .	24
3.2.1	Outcome generici vs specifici . . . . .	24
3.2.2	Trasformazione in utilità e calcolo dei QALY . . . . .	24
3.2.2.1	Utilità . . . . .	24
3.2.2.2	QALY . . . . .	24
<b>III</b>	<b>Introduzione</b>	<b>27</b>
<b>4</b>	<b>Studi biomedici</b>	<b>29</b>
4.1	Classificazioni . . . . .	29
4.2	Misurazione ed errori . . . . .	30

<b>5</b>	<b>Strumenti di aggiornamento e ricerca bibliografica</b>	<b>31</b>
5.1	Cosa sappiamo ad oggi . . . . .	31
5.1.1	Sintesi su malattie per professionisti . . . . .	31
5.1.2	Risposte più approfondite . . . . .	31
5.1.3	Ricerca bibliografica . . . . .	31
5.1.3.1	Dizionario . . . . .	31
5.1.3.2	Costruzione della ricerca . . . . .	31
5.1.3.3	Identificazione dei termini da ricercare . . . . .	33
5.2	Cosa bolle in pentola . . . . .	33
<b>6</b>	<b>Misure epidemiologiche assortite</b>	<b>35</b>
6.1	Misure e test di associazione . . . . .	35
6.1.1	Esposizione ed esito dicotomici . . . . .	35
6.1.1.1	Misure . . . . .	35
6.1.1.2	Test . . . . .	36
6.1.2	Esposizione multinomiale, esito dicotomico . . . . .	36
6.1.2.1	Misure . . . . .	36
6.1.2.2	Linear trend test . . . . .	36
6.1.2.3	Test di non linearità . . . . .	37
<b>7</b>	<b>Protocollo, raccolta dati e articolo</b>	<b>43</b>
7.1	Scrittura protocollo . . . . .	43
7.2	Dati e loro raccolta . . . . .	43
7.2.1	Tipologia di variabili . . . . .	43
7.3	Scrittura articolo . . . . .	44
7.3.1	Autorship . . . . .	44
<b>8</b>	<b>Confounding e interazione</b>	<b>47</b>
<b>IV</b>	<b>Studi sperimentali</b>	<b>53</b>
<b>9</b>	<b>Studi sperimentali</b>	<b>55</b>
9.1	Fasi degli studi sperimentali (farmacologici) . . . . .	55
<b>10</b>	<b>Fase 1</b>	<b>57</b>
10.1	Obiettivi . . . . .	57
10.2	Popolazione . . . . .	58
10.3	Definizioni . . . . .	58
10.4	Disegni . . . . .	58
10.4.1	Disegno standard (3+3) . . . . .	59
10.4.1.1	Funzionamento . . . . .	59
10.4.1.2	Critiche . . . . .	60
10.4.1.3	Varianti dello schema . . . . .	60
10.4.2	Continual reassessment method (Adattamento continuo) . . . . .	60
10.4.3	Disegno per nuovi farmaci . . . . .	61

<b>11 Fase 2</b>	<b>63</b>
11.1 Obiettivi . . . . .	63
11.2 Aspetti da considerare nel disegno . . . . .	63
11.2.1 Popolazione . . . . .	63
11.2.2 Trattamento . . . . .	64
11.2.3 Outcome . . . . .	64
11.2.4 Randomizzazione . . . . .	65
11.2.5 Scopo (sottofase) dello studio . . . . .	65
11.2.6 Categorie disegni . . . . .	65
11.3 Disegni comuni . . . . .	67
11.3.1 Livelli di errore nell'inferenza . . . . .	67
11.3.2 Stadio unico - A'Hern . . . . .	67
11.3.3 Due stadi - Simon . . . . .	68
11.3.4 Altri disegni . . . . .	70
11.3.5 Stima al termine di un multistage . . . . .	70
11.4 Criteri RECIST . . . . .	70
11.4.1 Classificazione delle lesioni e tumour burden . . . . .	70
11.4.2 Risposta . . . . .	71
11.4.3 Outcome derivabili . . . . .	72
<b>12 Feasibility/Pilot studies</b>	<b>73</b>
12.1 Definizioni . . . . .	73
12.2 Approccio a Maglietta - (mail mia 12/1/23) . . . . .	73
<b>13 Fase 3</b>	<b>75</b>
13.1 Obiettivi . . . . .	75
13.2 Classificazione di studi . . . . .	75
13.2.1 Studi esplicativi e pragmatici . . . . .	75
13.3 Validità di uno studio . . . . .	75
13.3.1 Validità interna . . . . .	76
13.3.2 Validità esterna . . . . .	76
13.4 PICO . . . . .	76
13.4.1 Popolazione . . . . .	76
13.4.1.1 Criteri di inclusione/esclusione . . . . .	76
13.4.1.2 Popolazione d'analisi . . . . .	77
13.4.2 Outcome . . . . .	77
13.5 Randomizzazione . . . . .	77
13.5.1 Alcuni concetti . . . . .	77
13.5.2 Tipologie . . . . .	78
13.5.3 Altre questioni . . . . .	79
13.6 Definizione dell'effetto del trattamento . . . . .	79
13.7 Disegni meno frequenti . . . . .	80
13.7.1 Disegno a più bracci paralleli . . . . .	80
13.7.2 Disegno fattoriale . . . . .	80
13.7.3 Disegno cross-over . . . . .	81
13.8 Altri studi comparativi (metodologicamente inferiori) . . . . .	81
13.8.1 Prima-dopo . . . . .	81
13.8.2 Trial controllati non randomizzati . . . . .	82
<b>14 Fase 4</b>	<b>83</b>

<b>V</b>	<b>Studi osservazionali</b>	<b>85</b>
<b>15</b>	<b>Introduzione osservazionali</b>	<b>87</b>
15.1	Tipi di studi . . . . .	87
15.2	Bias e confondimento . . . . .	87
<b>16</b>	<b>Coorte</b>	<b>89</b>
16.1	Disegni . . . . .	89
16.2	Pro/controllo . . . . .	90
16.3	Strategie d'analisi . . . . .	90
<b>17</b>	<b>Caso controllo</b>	<b>91</b>
<b>VI</b>	<b>Studi di diagnostica</b>	<b>93</b>
<b>18</b>	<b>Introduzione all diagnostica</b>	<b>95</b>
18.1	Introduzione . . . . .	95
18.1.1	Disegni di ricerca principali . . . . .	96
18.1.2	Architettura della ricerca diagnostica . . . . .	96
18.2	RCT Diagnostici . . . . .	96
<b>19</b>	<b>Studi di accuratezza diagnostica</b>	<b>97</b>
19.1	Introduzione . . . . .	97
19.2	Misure di accuratezza diagnostica . . . . .	98
19.2.1	Dati dicotomici . . . . .	98
19.2.1.1	Sensibilità, specificità . . . . .	98
19.2.1.2	Valori predittivi . . . . .	99
19.2.1.3	Uso di R - Stima accuratezza . . . . .	100
19.2.1.4	Uso di R - Inferenza accuratezza . . . . .	101
19.2.1.5	Molteplicità di focus diagnostici entro paziente . . . . .	102
19.2.2	Dati quantitativi . . . . .	102
19.2.3	Dati ordinali . . . . .	103
<b>VII</b>	<b>Revisioni sistematiche</b>	<b>105</b>
<b>20</b>	<b>Introduzione</b>	<b>107</b>
20.1	Definizioni e risorse utili . . . . .	107
20.2	Preparazione e mantenimento di una review (Cochrane) . . . . .	108
20.2.1	Protocollo . . . . .	108
20.2.2	Review team . . . . .	109
20.3	Domanda della ricerca e criteri inclusione . . . . .	109
20.4	Ricerca degli studi . . . . .	109
20.5	Selezione degli studi e collezione dati . . . . .	110
20.5.1	Selezione studi . . . . .	110
20.5.2	Dati da raccogliere . . . . .	111
20.5.2.1	Elementi . . . . .	111
20.5.2.2	Stime per misure dicotomiche . . . . .	112
20.5.2.3	Stime per variabili quantitative . . . . .	113
20.5.2.4	Stime per analisi di sopravvivenza . . . . .	113

20.6	Valutazione del rischio di bias negli studi inclusi . . . . .	113
20.6.1	Fonti di bias nei clinical trial . . . . .	113
20.7	Mantenimento della Revisione . . . . .	113
<b>21</b>	<b>Analisi dei dati e metanalisi</b>	<b>115</b>
21.1	Outcome e misure di efficacia . . . . .	115
21.2	Eterogeneità . . . . .	115
21.3	Metanalisi in a nutshell . . . . .	117
21.4	Metodi di calcolo dei pesi $W_i$ . . . . .	117
<b>22</b>	<b>Effect size and precision</b>	<b>119</b>
22.1	Overview . . . . .	119
22.2	Effect size basati su medie . . . . .	119
22.2.1	Differenza di medie non standardizzate in gruppi indipendenti . . . . .	119
22.2.2	Differenza di medie standardizzate in gruppi indipendenti	120
22.2.3	Response ratio . . . . .	121
22.3	Effect size basati su dati binari (tabelle $2 \times 2$ ) . . . . .	121
22.3.1	Risk ratio . . . . .	121
22.3.2	Odds ratio . . . . .	121
22.3.3	Risk difference . . . . .	122
22.3.4	Considerazioni sulla scelta . . . . .	122
22.4	Effect size basati su correlazioni . . . . .	122
22.5	Conversione tra effect size . . . . .	122
<b>23</b>	<b>Modelli ad effetti fissi e ad effetti random</b>	<b>123</b>
23.1	Introduzione . . . . .	123
23.2	Effetto fisso . . . . .	124
23.3	Effetti random . . . . .	124
23.4	Un confronto . . . . .	125
23.5	Esempi . . . . .	126
23.5.1	Dati dicotomici . . . . .	126
23.5.2	Dati continui . . . . .	130
23.5.3	Correlazioni . . . . .	131
<b>24</b>	<b>Eterogeneità</b>	<b>133</b>
24.1	Quantificazione . . . . .	133
24.1.1	Test di eterogeneità . . . . .	133
24.1.2	Scarto di eterogeneità . . . . .	134
24.1.3	Stima di $\tau^2$ . . . . .	134
24.1.4	$I^2$ . . . . .	134
24.1.5	Applicazioni in R . . . . .	134
24.2	Prediction intervals . . . . .	135
24.3	Analisi per sottogruppi . . . . .	135
24.4	Metaregressione . . . . .	135

<b>VIII</b>	<b>Dimensionamento campionario</b>	<b>139</b>
<b>25</b>	<b>Introduzione al dimensionamento campionario</b>	<b>141</b>
25.1	Approcci e ambiti di dimensionamento . . . . .	141
25.1.1	Errori nei test di ipotesi . . . . .	141
25.1.2	Giustificazione del dimensionamento . . . . .	141
25.1.3	Approcci al dimensionamento . . . . .	141
25.1.4	Ambiti di dimensionamento . . . . .	142
25.2	Ipotesi a confronto e disegni . . . . .	142
25.3	Considerazioni assortite . . . . .	143
25.3.1	Test a una o due code . . . . .	143
25.3.2	Aggiustamenti per dropouts . . . . .	143
25.3.3	Pacchetti R . . . . .	144
<b>26</b>	<b>Un gruppo</b>	<b>145</b>
26.1	Precision analysis - casi base . . . . .	145
26.1.1	Stima di una media . . . . .	145
26.1.1.1	Intervallo a due code . . . . .	145
26.1.1.2	Intervallo a una coda . . . . .	146
26.1.2	Stima di una proporzione . . . . .	147
26.1.2.1	Intervallo a due code . . . . .	147
26.2	Power analysis . . . . .	147
26.2.1	Test per una media . . . . .	147
26.2.1.1	Equivalenza . . . . .	148
26.2.1.2	Superiority/Non-inferiority . . . . .	151
26.2.1.3	Equivalence . . . . .	152
26.2.2	Test per una proporzione . . . . .	152
<b>27</b>	<b>Two groups</b>	<b>153</b>
27.1	T-test . . . . .	153
<b>28</b>	<b>Multiple endpoints</b>	<b>157</b>
28.1	Methodology . . . . .	157
28.2	MPE . . . . .	159



# Parte I

## Misc



# Capitolo 1

## Associazione tra variabili

### 1.1 Entrambe quantitative

### 1.2 Una quantitativa e una ordinale

### 1.3 Entrambe ordinali

#### 1.3.1 Kendall $\tau$ (tau)

Ipotizzando di aver rilevato due variabili X e Y su  $n$  individui, ottenendo  $(x_1, y_1), \dots, (x_n, y_n)$ ; per calcolarlo si costruiscono tutte le  $\binom{n}{2}$  coppie totali di osservazioni del dataset, e ciascuna di esse viene classificata come segue:

- la coppia composta da  $(x_i, y_i)$  e  $(x_j, y_j)$  è *concordante* se l'ordinamento delle due variabili coincide, ossia se alternativamente  $x_i > x_j$  e  $y_i > y_j$  oppure se  $x_i < x_j$  e  $y_i < y_j$  (in prima istanza si assume che non vi siano *ties*, ossia che non vi siano casi di  $x_i = x_j$  o  $y_i = y_j$ );
- la coppia è *discordante* in caso contrario.

##### 1.3.1.1 Prima versione: $\tau_A$

La prima versione di Tau si definisce come:

$$\begin{aligned}\tau_A &= \frac{\text{concordanti} - \text{discordanti}}{\text{totali}} \\ &= \frac{\sum_{i < j} \text{sgn}(x_i - x_j) \text{sgn}(y_i - y_j)}{\binom{n}{2}} \\ &= \frac{2}{n(n-1)} \sum_{i < j} \text{sgn}(x_i - x_j) \text{sgn}(y_i - y_j)\end{aligned}$$

Per l'interpretazione:

- se i ranghi delle due variabili coincidono il coefficiente ha valore 1
- se i ranghi delle due variabili sono uno l'inverso dell'altro il coefficiente ha valore -1

- se  $X$  e  $Y$  sono indipendenti, ci attendiamo che il coefficiente sia approssimativamente 0

### 1.3.1.2 Seconda versione: $\tau_B$

Qualora nel dataset siano presenti ties una data coppia può non essere classificata né concordante né discordante; Tau viene modificato lasciando invariato il numeratore e modificando il denominatore come segue:

$$\tau_B = \frac{\sum_{i < j} \text{sgn}(x_i - x_j) \text{sgn}(y_i - y_j)}{\sqrt{((\binom{n}{2} - n_1)(\binom{n}{2} - n_2))}}$$

con

- $n_1, n_2$  sono definiti come:

$$n_1 = \sum_i t_i(t_i - 1)/2$$

$$n_2 = \sum_j u_j(u_j - 1)/2$$

c

- $t_i$  il numero di valori *tied* nell' $i$ -esimo gruppo di ties per la prima variabile;
- $u_j$  il numero di valori *tied* nel  $j$ -esimo gruppo di ties per la seconda variabile.

### 1.3.1.3 Terza versione

Nel caso in cui le due variabili abbiano un numero di modalità differenti (per cui la crossclassificazione non è un quadrato ma un rettangolo) vi è una terza variante reputata migliore

$$\tau_C = \frac{2(n_c - n_d)}{n^2 \frac{m-1}{m}}$$

con

- $n_c$  numero di coppie concordanti
- $n_d$  numero di coppie discordanti
- $m = \min(r, c)$  con  $r$  numero righe (es modalità prima variabile) e  $c$  numero colonne (modalità seconda)

**L'utilizzo di R** `cor` specificando il metodo `kendall` produce  $\tau_A$  se non vi sono ties e  $\tau_B$  qualora ve ne siano. `cor.test` per il test ed intervallo di confidenza  $\tau_C$  non è implementata in R out of the box (poco male)<sup>1</sup>

---

<sup>1</sup>Per vedere come calcola la varianza SAS qui.

### 1.3.2 Goodman e Kruskal $\gamma$

È simile alla definizione di *tau*, dal quale differisce per il denominatore

$$\gamma = \frac{\text{concordanti} - \text{discordanti}}{\text{concordanti} + \text{discordanti}}$$

in sostanza ignorando i *ties* e focalizzandosi sulle coppie per le quali è possibile determinare un ordinamento. Dalla definizione si ha che  $\tau \leq \gamma$  (in valore assoluto), ossia che *tau* è generalmente più conservativo (in funzione di quanti *ties* ci sono).

variabili nominali Cramer V

## 1.4 Una ordinale e una nominale



Parte II

Health economics





## Capitolo 2

# Introduzione

### 2.1 Economia sanitaria

L'economia sanitaria applica la teoria economica all'analisi del mercato sanitario, e dei suoi problemi di allocazione.

La giustificazione dell'intervento pubblico nel campo sanitario può derivare da motivi efficientistici ed equitativi:

1. mancanza di sovranità del consumatore: situazione di informazione imperfetta ed asimmetrica, interventi sanitari come experience good
2. presenza di forti esternalità positive insite nel bene salute
3. la spesa per la salute è connessa ingran parte alla contrazione di malattie, quindi ad eventi aleatori: sarebbe pertanto efficiente l'impiego di contratti di assicurazione per la copertura dei rischi. La realizzazione di un sistema privato di assicurazione potrebbe però essere compromessa da probabilità di evento non indipendenti (es epidemie), vicine all'unità (es malattie ereditarie), problemi di moral hazard (es caso del terzo pagante) ed averse selection. Altro problema di carattere etico consiste nella possibilità di creare cream skinning da parte delle agenzie assicurative
4. mancanza di una molteplicità di produttori, soprattutto se si concentra l'attenzione in un contesto di piccole o medie dimensioni (provinciali/regionali)
5. disomogeneità del prodotto/prestazioni, anche se si considera la fornitura del medesimo servizio
6. forma di finanziamento delle prestazioni indiretta: fiscalità, non prezzi

Ovvio è (ma neanche troppo, da un punto di vista storico) che le decisioni allocative dello Stato debbano basarsi su studi comparativi delle alternative disponibili: ecco che entra in gioco la valutazione economica in sanità

### 2.2 Sistema sanitario e valutazione economica

Alla base della valutazione economica troviamo concetti classici dell'economia:

- **costo opportunità:** consiste nel valore dell'alternativa cui si rinuncia, impiegando le risorse in un dato modo. Nel campo sanitario esso è spesso rappresentato da outcome sanitari (quantificati in differenti unità di misura) raggiungibili mediante investimenti in alternative.  
L'utilizzo del costo opportunità è alla base della valutazione economica
- **costo marginale:** variazione del costo totale derivante dall'aumentare di una unità la produzione di un prodotto/servizio; questa è variabile fondamentale di orientamento della scelta di produzione perché in linea di massima, l'aumento del prodotto dovrebbe essere effettuato solo quando il beneficio marginale supera il costo marginale

Un **sistema sanitario è efficiente** (Dirindin) quando:

1. le risorse vengono impiegate per prestazioni efficaci, ossia in grado di incidere beneficamente sulla patologia del paziente
2. tali prestazioni presentano un beneficio marginale uguale o maggiore rispetto a quello ottenuto con impieghi alternativi
3. si inizia nella produzione da quelle erogazioni che garantiscono benefici di outcome maggiori, passando via via a quelle minori;
4. nella produzione del servizio occorre scegliere la combinazione di fattori produttivi che minimizzi i costi

Questo garantirebbe efficienza (ossia la maggior produzione di benessere possibile), ma non necessariamente equità (ossia distribuzione del benessere secondo criteri preferiti da un punto di vista etico-politico).

## 2.3 Processo di valutazione economica in sanità

Il perseguimento di un sistema sanitario efficiente parte necessariamente dalle singole azioni. ogni scelta allocativa (un farmaco/intervento piuttosto che un altro) dovrebbe fondarsi su un processo di valutazione articolato nelle seguenti fasi:

1. definizione della domanda di studio e scelta delle alternative confrontabili: occorre specificare chiaramente la popolazione/su che pazienti, il trattamento da valutare, quello con cui compararlo (es pratica standard), su che pazienti, e la variabile utilizzata per giudicarne l'effetto clinico.
2. analisi dei dati di efficacia della terapia: vi sono diversi studi e metodi di analisi di efficacia
3. **scelta tecnica di valutazione economica:** gli esiti sanitari debbono essere trasformati in output economici.  
Vi sono diversi approcci impiegabili che dipendono dalla disponibilità dei dati
  - nell'analisi costo efficacia i benefici del trattamento sono misurati in unità naturali; il problema è che studi che adottano un differente indicatore di efficacia non possono esser confrontati

- nell'analisi costo utilità si pesano gli anni di vita mediante la qualità di vita dei pazienti. Si utilizzano questionario che riportano su una scala normalizzata la condizione di salute. Ogni anno si riceve un punteggio compreso tra 0 (per la morte) e 1 (in caso di perfetta salute); questi punteggi vengono detti QOL (quality of life): ad esempio un individuo di media salute (es  $qol=0.5$ ) necessita di due anni di vita per guadagnarsi un QALY (quality adjusted life year). Mediante un confronto tra l'esito di tale indicatore per i trattati e per il gruppo di controllo è possibile ottenere una stima della variazione della qualità della vita indotta dal trattamento
- nell'analisi costi benefici, si esprime il beneficio clinico in termini monetari, motivo per cui questo approccio è stato criticato

Per quanto riguarda i costi essi possono essere classificati in

- diretti: costi associati alla produzione del servizio sanitario (ricoveri, terapie, visite mediche) e non (assistenza domiciliare)
  - indiretti: perdite di produttività del paziente e dei rispettivi familiari
4. **scelta ottica di analisi valutativa** (es sistema sanitario pagante, società): ha impatto sui costi e benefici che si dovranno considerare. Ad esempio se fatto nell'ottica del sistema sanitario dobbiamo tenere conto degli stimendi dei medici, non delle perdite di produttività di pazienti e familiari.
5. **svolgimento dell'analisi economica**: attualizzazione, calcolo di indicatori sintetici, e analisi di sensibilità. Si ha

- attualizzazione di costi e benefici avviene con un tasso compreso tra il 3 e il 5
- sia effetto che costo possono essere maggiori, minori o uguali rispetto all'alternativa comparata.  
La scelta allocativa non appare evidente se si hanno contemporaneamente maggiori benefici e maggiori costi, oppure minori benefici e minori costi.  
In questi casi (per analisi costo efficacia e costo utilità) occorre rapportare l'incremento dei costi con quello dei benefici. Se si è utilizzato un approccio costo-efficacia sarà
- analisi di sensibilità: ve ne sono due tipi
  - (a) si fanno variare i parametri utilizzati nel computo di costi e benefici (es tasso di interesse, funzioni di distribuzione di costi e benefici) per vedere come cambia il risultato finale. Facendo variare un parametro alla volta si ha una analisi univariata, più parametri contemporaneamente. I risultati sono robusti se le conclusioni non variano
  - (b) si fanno best e worst case scenario: best case si considera efficacia maggiore e costi minori, worst case il contrario

## 2.4 Cost utility shit

Qol possono esser calcolati secondo due approcci:

- approccio *clinico*: si valuta la condizione di salute del paziente mediante questionari medici generici (SF36, EQ5D) o specifici per patologia. Gli strumenti generici hanno dalla loro il beneficio di versatilità d'uso, quelli specifici dell'approfondimento di indagine
- approccio *economico*: si chiede al paziente di valutare, in maniera diretta o meno, il proprio stato di salute (approccio che si richiama maggiormente al concetto di utilità. Vi sono tre tecniche principali:
  1. **VAS** o rating scale: si richiede al paziente di porre una X su una linea che va dalla perfetta salute alla morte per indicare quanto ci si ritiene in salute
  2. **standard gamble**: si chiede al paziente di scegliere tra un determinato stato di malattia e la partecipazione all'urna che da  $p$  probabilità di morire e  $(1 - p)$  di vivere in perfetta salute; qual è la  $p$  che rende indifferente la scelta.
  3. **time trade off**: viene chiesto al paziente quanti anni di vita nello stato di malattia sarebbe disposto a rinunciare per vivere in salute. L'utilità è pari a

$$u = \frac{\text{anni di vita in condizioni perfette se si accetta}}{\text{anni di vita da vivere in malattia}} \quad (2.1)$$

Idea è integrare la valutazione di efficacia con quella dei costi, al fine di dire se qualcosa di sperimentale è il miglior risultato di salute ottenibile con le risorse a disposizione.

**Prospettive di analisi** Si possono adottare le prospettive del:

- paziente (e famiglia),
- del terzo pagante (SSN): qui non sono rilevanti i costi privati di pazienti e familiari
- più in generale della società:

**Valutazione dei costi** I costi di un programma/trattamento etc possono essere diretti, indiretti o intangibili, e sono misurati in unità monetarie:

- i costi diretti sono quelli imputabili alla malattia principale (+ complicanze da co-morbosità) del settore sanitario (personale, attrezzature, farmaci, materiali), di pazienti e familiari (trasporto, alloggio, assistenza domiciliare) e di altri settori (modifiche richieste dal programma)
- i costi indiretti (o sociali) sono quelli del tempo di lavoro perso (paziente, familiari) e riduzione della produttività
- i costi intangibili (es dolore, ansia per intervento) possono solo essere elencati (sono compresi nei QALY)

**Valutazione delle conseguenze** Le conseguenze possono essere misurate in differenti modi, dando origine a diversi tipi di analisi, quando si va a confrontare coi costi/benefici tra due o + alternative:

- unità fisiche (es decessi, infezioni) da origine a costo-efficacia
- utilità (QALY) e costo-utilità
- unità monetarie (benefici diretti, indiretti e intangibili) da origine ad analisi costi-benefici

## 2.5 Comparazione di costi ed effetti

**Piano costi-efficacia** La differenza tra costi ed esiti nei due o più programmi messi a confronto viene poi plottata su un diagramma cartesiano con asse  $x$  la differenza di efficacia e asse  $y$  la differenza di costi:

- nel quadrante sud-est il trattamento sperimentale ha esiti migliori e costi minori, pertanto è palesemente superiore al trattamento controllo
- nel quadrante nord-ovest il trattamento sperimentale ha efficacia minore e costo maggiore ed è palesemente dominato dal trattamento controllo
- nei quadranti nord-est e sud-ovest non vi è un trattamento manifestamente superiore ed è necessaria una valutazione ulteriore: a nord est costi ed effetti del trattamento sperimentale sono superiori, mentre in sud-ovest si ha un effetto minore ma anche un costo. Bisognerà in entrambi i casi confrontare il costo per unità di efficacia con una soglia, che costituisce la pendenza di una retta passante dal centro degli assi. Le soluzioni al di sotto della linea sono considerate accettabili.

**ICER** A questo punto è possibile calcolare l'ICER (incremental cost effectiveness ratio) come

$$ICER = \frac{\Delta C}{\Delta E} = \frac{C_{exp} - C_{control}}{E_{exp} - E_{control}}$$

che dice il costo aggiuntivo per unità di effetto aggiuntivo. Le analisi di costo efficacia e costo utilità differiranno per il denominatore, dove avremo rispettivamente unità naturali (anni vita) o QALY (anni vita moltiplicati per QOL)

Il programma è considerato finanziabile se ha un ICER minore della disponibilità a pagare per unità di outcome aggiuntivo

$$\frac{\Delta C}{\Delta E} < R_T \quad (2.2)$$

con  $R_t$  la soglia che manifesta la disponibilità a pagare per una unità di efficacia in più (es anno di vita/QALY in più).

**Aggiunta della variabilità** La stima dell'ICER ad ora è unica e non dispone di variabilità. Il modo più comune per aggiungerla è mediante bootstrap: mediante questo si avranno

- una nuvola di punti sul piano costo efficacia
- una stima dell'intervallo di confidenza dell'ICER (es BCA bootstrap)

Per valori differenti di soglia si può calcolare la percentuale di campioni che rispettano la 2.2, detta probabilità di costo efficacia.

**NMB** L'equazione 2.2 di scelta dell'ICER si può riarrangiare portando tutto a destra e riesprimendo il criterio di accettazione sulla base del Net Monetary Benefit come:

$$NMB = R_t \Delta E - \Delta C > 0$$

In questa formulazione il programma è accettabile rispetto alla controparte se ha un NMB positivo. Dato che queste valutazioni dipendono fortemente da  $R_t$  ossia la disponibilità a pagare del finanziatore si può calcolare:

- il NMB per vari livelli di soglia, per vedere dove diventa positivo
- nel caso in cui vi sia variabilità nelle stime di costi e ricavi, al variare di  $R_t$ , quale è la probabilità che l' $NMB > 0$

**Regressioni: OLS e SUR** Dato che ogni paziente ha una misurazione di efficacia, una di costo e si ha una disponibilità a pagare, la NMB può essere calcolata a livello individuale come

$$NMB_i = E_i \cdot R_t - C_i$$

questo apre spazio alla modellazione mediante OLS, es per quantificare l'effetto del trattamento sul NMB medio più aggiustamento per covariate varie, in setting osservazionali o anche sperimentali.

Metodi più recenti (SUR) si sono spostati sul definire due equazioni di regressione, una per costi e una per gli effetti e correlandone il termine d'errore (non considerato indipendente tra i due: tipicamente costi ed efficacia sono negativamente correlati, costi maggiori si hanno in condizioni di efficacia minore). L'approccio seemingly unrelated regression ha il vantaggio di permettere differenti covariate e forme funzionali nelle due equazioni. Vedere Willan Briggs e Hoch 2004: regression methods for covariate adjustment and subgroup analysis for non censored cost-effectiveness data

## Capitolo 3

# Outcomes

### 3.1 Costi

Il metodo di raccolta dati sui costi dipende se siamo all'interno di un trial, dove accesso/quantificazione a livello di singolo paziente è possibile oppure no

- nei trial un modo per quantificare le prestazioni è il questionario CSRI (Client Services Receipt Inventory)
- al di fuori ci sono valorizzazioni di DRG

Più info sulla valorizzazione dei costi in [17]

#### 3.1.1 Attualizzazione

Costi che si verificano in diversi periodi debbono essere attualizzati, scontandoli ad un determinato tasso ([5] riporta che NICE consiglia il 3.5% ma meglio effettuare analisi di sensibilità tra lo 0 e il 6).

```
## esempio pag 18 baio
npv <- function(y, t, i){
  #y flow
  #t time of flow in years
  #i yearly interest rate
  num <- y
  den <- (1 + i)^t
  sum(num/den)
}

npv(rep(15000, 5), 0:4, 0.035)

## [1] 70096.19
```

## 3.2 Efficacia

### 3.2.1 Outcome generici vs specifici

La qualità della vita del paziente deve essere tradotta in una utilità, ossia una misura nell'intervallo 0-1 con 0 = morte e 1 = perfetta salute.

Per misurare la qualità della vita si possono utilizzare strumenti generici o specifici:

- decision maker tendono a preferire outcome generici, che permettono di comparare costo-efficacia tra varie malattie/aree;
- laddove vi sia preoccupazione che uno strumento generico possa non funzionare bene o non vi siano dati di comparazione si può ricorrere a uno strumento specifico, posto che sia

Tra gli strumenti generici si annoverano:

- SF-6D: 11 item presi dall'SF-36 e combinati per produrre uno score
- EQ-5D-5L: scala in due parti (la prima su aspetti specifici, la seconda una valutazione overall)
- FACT-8D

### 3.2.2 Trasformazione in utilità e calcolo dei QALY

#### 3.2.2.1 Utilità

Approfondire la metodologia in R nei pacchetti `eq5d` o `fact-8D` (per il FACT-G).

#### 3.2.2.2 QALY

Dopo il calcolo dell'utilità a diversi punti occorre calcolare i QALY: l'attesa di vita di 1 anno di salute perfetta è valorizzata 1, mentre morte è valorizzata con 0.

Occorre che vi sia una valutazione di benessere anche al baseline. Vedere comunque <https://github.com/Health-Economics-in-R/QALY> e [https://en.wikipedia.org/wiki/Quality-adjusted\\_life\\_year](https://en.wikipedia.org/wiki/Quality-adjusted_life_year) per altre robe sui QALY

```
QALY <- function(u, t, i = 0){
  ## browser()
  ## u utility
  ## t times in years
  ## i interest rate
  if (t[1] != 0) stop("QALY needs baseline measurement for utility")
  durate <- diff(t)
  utility_couples <- cbind(u[-length(u)], u[-1])
  mean_u <- apply(utility_couples, 1, mean, na.rm = TRUE)
  sum(durate * mean_u * (1/(1 + i)^(t[-length(t)])))
}
```

*##y flow*



```

##t time of flow in years
##i yearly interest rate

## num <- y
## den <- (1 + i)^t

## baio pag 24
u <- c(0.656, 0.744, 0.85, 0.744, 0.744)
t <- c(0, 6, 12, 18, 24)/12
QALY(u, t)

## [1] 1.519

u <- c(0.656, 0.656, 0.656, 0.656, 0.744)
QALY(u, t)

## [1] 1.334

```

**Attualizzare i QALY** Diversi attualizzano il QALY applicando un tasso di sconto (comune ai costi).

```

## variazione con attualizzazione
QALY(u, t, i = 0.035)

## [1] 1.299712

```



## Parte III

# Introduzione



# Capitolo 4

## Studi biomedici

### 4.1 Classificazioni

**Definition 4.1.1** (Oggetto dello studio). Si hanno

- studi per *descrivere una casistica* in un dato momento (ad esempio osservazionali cross-section)
- studi ove si cercano *relazioni di causa effetto* (efficacia trattamento, studi prognostici, studi eziologici)
- studi *diagnostici*
- studi per lo *sviluppo di strumenti/tool* (ad esempio questionari)
- studi che *sintetizzino l'evidenza disponibile* (revisioni sistematiche e meta-analisi)

**Definition 4.1.2** (Ruolo del ricercatore). In studi dove si indagano relazioni di causa-effetto possiamo distinguere[4]:

- **Osservazionali**: il ricercatore studia la relazione tra una *caratteristica* (fattore demografico, ambientale, marker genetico), ed una *variabile di interesse* (detta anche *outcome*, ad esempio insorgenza malattia o guarigione), senza intervenire in alcun modo sulle condizioni in cui lo studio viene condotto: seleziona il campione e poi osserva eventuali associazioni tra possibili fattori di rischio/protezione e la malattia che possano *suggerire* una relazione di causa-effetto.
- **Sperimentali** studi nei quali il ricercatore controlla le condizioni di svolgimento; nello specifico il ricercatore studia la relazione tra un *fattore sperimentale*, assegnato dal ricercatore stesso, ed un *variabile di interesse*; tutto questo il più possibile al netto dell'effetto di *fattori sub-sperimentali* (caratteristiche demografiche o della malattia, trattamenti precedenti e concomitanti, centro) che possano influenzare autonomamente la variabile di interesse.

**Definition 4.1.3** (Criteri causalità). Quando vi è causalità tra una caratteristica/fattore sperimentale e una variabile di interesse? alcune condizioni (nessuna necessaria o sufficiente) che rafforzano l'evidenza di causalità [25]

1. forte associazione tra caratteristica/fattore sperimentale e una variabile di interesse
2. l'esposizione a caratteristica/fattore sperimentale è temporalmente precedente alla manifestazione della variabile di interesse
3. vi è una plausibile spiegazione biologica
4. l'associazione è supportata da altre ricerche in setting differenti
5. reversibilità: vi dovrebbe essere evidenza che se la causa è rimossa, anche l'effetto dovrebbe scomparire
6. dose-effetto: vi dovrebbe essere evidenza che a maggiori livelli di esposizione maggiore è il manifestarsi della variabile di interesse
7. assenza di spiegazioni alternative convincenti per eventuali variazioni della variabile di interesse

## 4.2 Misurazione ed errori

*Remark 1.* Una misura può essere vista come la somma di diverse componenti

$$\text{valore misurato} = \text{valore reale} + \text{errore casuale} + \text{errore sistematico}$$

Il processo di misurazione di un qualsiasi fenomeno può essere riguardato da due tipi di errori e questi sono parte integrante della variabilità dei fenomeni come li conosciamo.

**Definition 4.2.1** (Errore casuale). Errore che produce oscillazioni che non hanno un andamento riproducibile

**Definition 4.2.2** (Errore sistematico (*bias*)). Errore che produce risultati che differiscono dal valore vero sistematicamente sempre nella stessa direzione

*Remark 2.* Sia errori casuali che sistematici hanno un impatto sul risultato della misurazione tuttavia il loro effetto è diverso:

- gli errori casuali, all'aumentare del numero delle misurazioni, tendono ad avere un impatto minore sugli indici di tendenza centrale (media/mediana) perché gli errori in una direzione o nell'altra tendono a compensarsi
- i secondi invece non si mitigano all'aumentare del campione

## Capitolo 5

# Strumenti di aggiornamento e ricerca bibliografica

### 5.1 Cosa sappiamo ad oggi

#### 5.1.1 Sintesi su malattie per professionisti

Per una sintesi su una malattia fare un mix delle seguenti:

- dynamed
- ClinicalKey; conviene selezionare **Clinical Overview** se si vuole una sintesi su una malattia
- UpToDate (sono sintesi fatte da esperti)

#### 5.1.2 Risposte più approfondite

Il database da consultare (soprattutto per diagnosi e trattamento) è quello delle <https://www.cochranelibrary.com/> che costituiscono il gold standard

#### 5.1.3 Ricerca bibliografica

##### 5.1.3.1 Dizionario

**Definition 5.1.1** (MEDLINE). Un database di studi contenente autori, abstract e link e indicizzati secondo MeSH

**Definition 5.1.2** (PubMed). Un motore di ricerca gratuito su MEDLINE

**Definition 5.1.3** (MeSH). Vocabolario di temi medici gerarchico utilizzato per indicizzare per topic gli articoli e permettere ricerca precisa

##### 5.1.3.2 Costruzione della ricerca

Prima di effettuare una ricerca nel database PubMed preprocessa la stringa di ricerca; per controllare cosa effettivamente PubMed cerchi (in seguito a quanto inserito) andare a vedere la history sotto l'advanced search.

Quando si effettua una ricerca semplice in Pubmed di default ci vengono restituiti molti risultati perchè:

- viene ricercato questo termine in tutti i campi del database (quindi nel titolo, autore, affiliazioni, giornale, abstract, lingua, mesh, keyword fornite dagli autori degli articoli). Questo è il modo più inclusivo di effettuare ricerche ma il rischio di falsi positivi è alto
- se il termine non è posto tra apici il termine viene applicato il mapping ossia si ammettono sinonimi e parole simili

Se vogliamo essere più precisi possiamo:

- specificare i fields
- disattivare il mapping (magari una volta individuati mesh e fields adatti) ponendo tra apici
- combinare il tutto con gli operatori logici e parentesi tonde

**Definition 5.1.4** (Operatori logici). Vanno categoricamente maiuscoli e sono

- AND: recupera items che contengono entrambi i termini
- OR: recupera items che contengono almeno uno dei termini
- NOT: recupera documenti che contengono solo il primo dei due termini, escludendo il secondo o i documenti in cui ci sia compresenza dei due

**Definition 5.1.5** (Parentesi tonde). Indirizza la precedenza/sequenza della ricerca come in algebra, ad esempio

`cancer AND (prognosis OR diagnosis)`

cercherà prima le parole prognosis o diagnosis e in seguito matcherà con i risultati che hanno anche cancer

**Definition 5.1.6** (Field tags). Si possono rinvenire nella ricerca avanzata, ma i più importanti sono:

- [au]: autore
- [mh] o [mesh]: mesh
- [ti]: cerca nel titolo. Utile per restringere di molto la ricerca
- [tiab]: cerca la parola nel titolo o nell'abstract
- [tw]: per essere più inclusivi rispetto a tiab cerca in titolo o abstract o mesh o nelle keyword fornite dall'autore più altro (ma non cerca in affiliazione, autore, giornale)

*Remark 3.* Per le revisioni sistematiche i più comuni sono tiab e tw

*Remark 4.* Il field tags si applicano a tutte le parole che precedono; nel seguente esempio tiab è applicato anche a cancer

`cancer prognosis[tiab] survival`

viene mappato a

`"cancer prognosis"[Title/Abstract] AND ("mortality"[MeSH Subheading] OR "mortality"[A`

**Definition 5.1.7** (Troncamento). L'utilizzo dell'asterisco funziona come carattere di globbing: es `osteo*` trova osteosarcoma, osteoarthritis etc



### 5.1.3.3 Identificazione dei termini da ricercare

Per ricercare con i termini corretti da utilizzare nella ricerca “definitiva”:

- analizzare i risultati di una ricerca base (in advanced -> history)
- analizzare i dati/categorizzazioni/metadati di un articolo considerato rilevante
- utilizzare il database delle MeSH (dalla schermata iniziale andare su Explore -> Mesh Database): per farlo cercare parole e poi aggiungere i filtri per la mesh che meglio rappresentano

## 5.2 Cosa bolle in pentola

Protocolli:

- ClinicalTrials.gov per singoli studi (trial ma anche osservazionali)
- prospero per le revisioni sistematiche
- cercare protocol nel titolo su pubmed



## Capitolo 6

# Misure epidemiologiche assortite

### 6.1 Misure e test di associazione

#### 6.1.1 Esposizione ed esito dicotomici

Nel caso sia l'esposizione (a un fattore di rischio o a un trattamento) che l'esito siano variabili dicotomiche le frequenze possono essere rappresentate in una tabella  $2 \times 2$  analoga a 6.1.

##### 6.1.1.1 Misure

Sulla base di essa sono definibili diverse misure di associazione tra esposizione ed esito.

**Definition 6.1.1** (Rischio). Definito come:

$$Risk = \frac{a+b}{n}$$

**Definition 6.1.2** (Risk ratio (o relative risk)). Definito come:

$$RR = \frac{a/(a+c)}{b/(b+d)} = \frac{a(b+d)}{b(a+c)}$$

con al numeratore il rischio tra gli esposti, al denominatore quello dei non esposti.

	Esposti	Non Esposti	Tot
Malati	$a$	$b$	$a+b$
Non Malati	$c$	$d$	$c+d$
Tot	$a+c$	$b+d$	$n = a+b+c+d$

Tabella 6.1: Tabella per misure di associazione

**Definition 6.1.3** (Risk difference (o absolute risk reduction)). Definita come:

$$RD = \frac{a}{a+c} - \frac{b}{b+d}$$

con al numeratore il rischio tra gli esposti, al denominatore quello dei non esposti.

**Definition 6.1.4** (Odds ratio). Definito come

$$OR = \frac{a/c}{b/d} = \frac{ad}{bc}$$

con al numeratore l'odd tra gli esposti, al denominatore quello dei non esposti.

*Remark 5.* Costituisce una buona approssimazione al rischio relativo quando si è in presenza di una malattia con prevalenza rara. Ossia:

$$a+c \simeq c, b+d \simeq d \iff RR = \frac{a/(a+c)}{b/(b+d)} \simeq \frac{a/c}{b/d} = OR \quad (6.1)$$

**TODO:** relazione odds ratio e risk ratio

**Relative risk e odds ratio** Generalmente presentare i risultati in termini di rischio relativo può essere compreso maggiormente (essendo i rischi delle probabilità). Le argomentazioni a favore dell'odds ratio:

- in alcuni casi (es studi caso controllo) è l'unica misura che può essere stimata correttamente;
- più facilmente stimabile (regressione logistica).

#### 6.1.1.2 Test

Per testare l'associazione tra i due fattori `fisher.test` o `chisq.test`, oppure costruire l'intervallo di confidenza degli stimatori e verificare che non includano la soglia di non differenza (per RD è 0, per OR e RR è 1).

### 6.1.2 Esposizione multinomiale, esito dicotomico

#### 6.1.2.1 Misure

L'esposizione multinomiale può essere naturale o derivare da una discretizzazione di una variabile quantitativa.

Comunque sia, nel caso si sceglie un gruppo di esposizione a fungere da gruppo base e si comparano gli altri livelli con questo utilizzando le misure già definite.

#### 6.1.2.2 Linear trend test

Un test che assume importanza in questo ambito è il test dei trend lineari [3], dove si verifica che a livelli via via crescenti del fattore la percentuale si modifichi linearmente (in aumento o decremento):

```

# x e n sono le frequenze nell'ordine del fattore di rischio crescente
x <- c(83, 90, 129, 70)
n <- c(86, 93, 136, 82)

## implementato in R con
prop.trend.test(x, n)

##
## Chi-squared Test for Trend in Proportions
##
## data:  x out of n ,
## using scores: 1 2 3 4
## X-squared = 8.2249, df = 1, p-value = 0.004132

## under the hood è un modello lineare sulle probabilità con pesi particolari
## dei quali non ho capito la genesi/interpretazione per ora
score <- 1:4
p <- sum(x) / sum(n)
w <- n/p/(1 - p)
df <- data.frame(freq = x/n, score = score)
a <- anova(mod <- lm(freq ~ score, data = df, weights = w))
chisq <- c(`X-squared` = a["score", "Sum Sq"])
(p <- pchisq(as.numeric(chisq), 1, lower.tail = FALSE))

## [1] 0.004131897

## riproduzione test su table 2.2 (pag 41) effettuato a pag 145 woodward
## x e n sono le frequenze nell'ordine del fattore di rischio crescente
x <- c(100, 382, 183, 668, 279, 109)
n <- c(592, 2254, 1017, 3150, 1253, 415)
prop.trend.test(x, n)$statistic

## X-squared
## 33.63156

```

Qualora il test venga significativo mentre il chi quadrato non lo sia sta a significare che sebbene non vi sia evidenza che le proporzioni negli strati differiscono dalla proporzione media, all'aumentare dello strato di esposizione si nota un incremento della proporzione registrata.

### 6.1.2.3 Test di non linearità

Il test si calcola la differenza tra l'astatistica chi quadrato e la statistica del trend test; tale differenza è comparata ai valori critici della ristribuzione chi quadrato con gradi di libertà dati dalla differenza delle due componenti

```

prop.nonlinear.trend.test <- function(x, score = seq_len(ncol(x))) {
  ## x is a 2 x l matrix (l is the number of groups
  ## todo farla anche per table
  ## e per due variabili (forse?)
  if (!is.matrix(x) || nrow(x) != 2) stop('x must be a 2 x l matrix')

```

```

method <- 'Test for non-linear trend in Proportions'
dname <- paste(paste(x[1,], collapse = ', '), "out of",
               paste(colSums(x), collapse = ', '))
dname <- paste(dname, "\n using scores:", paste(score, collapse = " "))
chi <- chisq.test(x)
ltt <- prop.trend.test(x = x[1, ], n = colSums(x))
test <- chi$statistic - ltt$statistic
df <- chi$parameter['df'] - ltt$parameter['df']
p_value <- pchisq(q = test, df = df, lower.tail = FALSE)
list(chi, ltt, test, df, p_value)
structure(list(statistic = test,
               data.name = dname,
               parameter = c("df" = df),
               p.value = p_value,
               method = method),
          class = 'htest')
}

## woodward pag 146

## table 2.2 pag 41
chd <- c(100, 382, 183, 668, 279, 109)
n <- c(592, 2254, 1017, 3150, 1253, 415)
nochd <- n - chd
m <- rbind(chd, nochd)
colnames(m) <- c("I", "II", "IIIIn", "IIIIm", "IV", "V")
m

##           I    II IIIIn IIIIm IV    V
## chd    100   382   183   668 279 109
## nochd 492 1872   834 2482 974 306

## validazione esempio pag 146
prop.nonlinear.trend.test(m)

##
## Test for non-linear trend in Proportions
##
## data: 100, 382, 183, 668, 279, 109 out of 592, 2254, 1017, 3150, 1253, 415
## using scores: 1 2 3 4 5 6
## X-squared = 2.7677, df.df = 4, p-value = 0.5974

prop.trend.test(m[1,], colSums(m))

##
## Chi-squared Test for Trend in Proportions
##
## data: m[1, ] out of colSums(m) ,
## using scores: 1 2 3 4 5 6
## X-squared = 33.632, df = 1, p-value = 6.66e-09

```

```
## -----
## attributable risk (woodward pag 148)
## -----
library(attribrisk)

## Error in library(attribrisk): non c'è alcun pacchetto chiamato
## 'attribrisk'

es3.1 <- matrix(c(31, 1386, 15, 1883), ncol = 2)
dimnames(es3.1) <- list('death' = c('Yes', 'No'), 'smoke' = c('Yes', 'No'))
## calcolo a mano (facendo riferimento a therneau come notazione di formule)
addmargins(es3.1)

##      smoke
## death Yes  No  Sum
##   Yes   31   15   46
##   No  1386 1883 3269
##   Sum 1417 1898 3315

pd <- 46/3315
pd_given_notf <- 15/1898
(ar <- (pd - pd_given_notf)/pd)

## [1] 0.4304646

## simuliamo una situazione di dataset reale
library(lbmisc)
es3.1_df <- table2df(as.table(es3.1))
## e' importante che i si e i no siano codificati bene (No gruppo base, Yessa
## l'evento per cui si vuole stimare il rischio attribuibile)
es3.1_df$death <- relevel(es3.1_df$death, ref = 'No')
es3.1_df$smoke <- relevel(es3.1_df$smoke, ref = 'No')
str(es3.1_df)

## 'data.frame': 3315 obs. of 2 variables:
## $ death: Factor w/ 2 levels "No","Yes": 2 2 2 2 2 2 2 2 2 ...
## $ smoke: Factor w/ 2 levels "No","Yes": 2 2 2 2 2 2 2 2 2 ...

table(es3.1_df)

##      smoke
## death  No  Yes
##   No 1883 1386
##   Yes   15   31

## table(es3.1_df)
attribrisk(death ~ expos(smoke), data = es3.1_df)

## Error in attribrisk(death ~ expos(smoke), data = es3.1_df): non
## trovo la funzione "attribrisk"
```

```

## non vengono proprio uguali gli intervalli di confidenza (il
## pacchetto usa jackknife, non una formula chiusa), comunque la stima
## c'e

es_th <- matrix(c(938, 763, 384, 559), ncol = 2)
dimnames(es_th) <- list('stroke' = c('Yes', 'No'), 'hbp' = c('Yes', 'No'))
addmargins(es_th)

##          hbp
## stroke Yes No Sum
##    Yes  938 384 1322
##    No   763 559 1322
##    Sum 1701 943 2644

es_th <- table2df(as.table(es_th))
es_th$stroke <- relevel(es_th$stroke, ref = 'No')
es_th$hbp <- relevel(es_th$hbp, ref = 'No')
(example1 <- attribrisk(death ~ expos(hbp), data = es_th))

## Error in attribrisk(death ~ expos(hbp), data = es_th): non trovo
## la funzione "attribrisk"

## COMUNQUE FARSI UNA FUNZIONCINA BASE BASE PER CALCOLARLO A MANO CHE
## NON MI FIDO ECCESSIVAMENTE

## -----
## rate
## -----

## esesmpio woodward pag 155

## rate uomini (1080/612955) (per 1000 uomini)
do.call(c, poisson.test(x = 1080, T = 612955)[c("estimate", 'conf.int')])*1000

## estimate.event rate          conf.int1          conf.int2
##           1.761956           1.658427           1.870256

## relative rate uomini vs donne
do.call(c, poisson.test(x = c(1080, 306),
                        T = c(612955, 634103))[c("estimate", 'conf.int')])

## estimate.rate ratio          conf.int1          conf.int2
##           3.651183           3.212988           4.158972

## stessa cosa fatta con modello di poisson
retrate_test <- data.frame(ev = c(1080, 306),
                          group = c(1, 0),
                          midp = c(612955, 634103))
rate_mod <- glm(ev ~ group + offset(log(midp)),
               data = retrate_test,
               family = poisson)
exp(cbind('estimate' = coef(rate_mod), confint(rate_mod)))['group', ]

```



```
## Waiting for profiling to be done...  
  
## estimate      2.5 %    97.5 %  
## 3.651183 3.220635 4.151744  
  
## il denominatore del rate in questo caso è la popolazione a metà  
## anno ma possono anche essere il tempo complessivo di osservazione
```



## Capitolo 7

# Protocollo, raccolta dati e articolo

### 7.1 Scrittura protocollo

Idee:

- partire dalla fine: ossia dalla pubblicazione che si vuole sottomettere e in altre parole che risposta si vuol dare, suoi contenuti e tipo di studio
- trovare una guideline per quel tipo di pubblicazione (consort, strobe, stard ecc) nell'equator network
- iniziare a compilare *direttamente in inglese*, il protocollo aiutandosi (ossia dando risposta a) tutti/maggior parte degli elementi messi in evidenza dalla guideline di pubblicazione relativa

Scrivere il protocollo direttamente in inglese fa sì che al momento della scrittura dell'articolo vero e proprio sezioni come Introduzione, Materiale e Metodi siano già state scritte e possano essere solo copiate e incollate; mancheranno solo i risultati derivanti dall'analisi statistica e la discussion.

Una volta scritto il protocollo e individuati i task da svolgere, individuare la lista degli autori della pubblicazione.

### 7.2 Dati e loro raccolta

#### 7.2.1 Tipologia di variabili

- **qualitative**: le rispettive modalità sono rappresentate da *attributi* (generalmente parole o frasi)
  - *nominali* (o *sconnesse*): le modalità non assumono alcun ordine. L'unica operazione effettuabile tra due unità consiste nello stabilire se posseggono o meno la stessa modalità/attributo.

**Example 7.2.1.** GENERE

- *ordinali*: variabile le cui modalità presentano un ordinamento, per cui è possibile stabilire tra le modalità di due unità una relazione di ordine rispetto alla variabile considerata.

**Example 7.2.2.** Per il TITOLO DI STUDIO,  $Licenza\ media < Diploma < Laurea$

- **quantitative**: le rispettive modalità sono rappresentate da numeri
  - *discrete*: le sue modalità (numero) possono essere poste in corrispondenza con l'insieme  $\mathbb{N}$  o un suo sottoinsieme (ossia  $\mathcal{Y}$  è finito o numerabile)

**Example 7.2.3.** NUMERO DI FIGLI

- *continue*: le sue modalità (numero) possono essere poste in corrispondenza con l'insieme  $\mathbb{R}$  o un suo sottoinsieme.

**Example 7.2.4.** ALTEZZA

*Remark 6.* Alcune variabili concettualmente continue (es età, altezza) possono essere registrate mediante valori discreti a causa della limitatezza di precisione insita nel relativo strumento di misurazione.

Le variabili quantitative possono anche essere classificate in base alla presenza di uno zero convenzionale o meno:

- *quantitativa per intervallo*: variabili che hanno una unità di misura ma non dello 0 (inteso come assenza della quantità da misurare), che viene inteso come convenzionale/arbitrario.  
Variabili del genere permettono solo un confronto per differenza tra le modalità che i soggetti assumono, mentre non ci permettono di calcolare rapporti che abbiano un senso (perché lo 0 della scala è arbitrario).  
Esempi: misurazioni della temperatura in Celsius o Fahrenheit; il tempo misurato su differenti calendari.
- *quantitativa per rapporto*: variabili per le quali è intrinseca/univoca la definizione dello zero, corrispondente all'assenza della caratteristica misurata. Valori negativi non dovrebbero esser possibili.  
Il fatto che l'origine sia condivisa permette di calcolare rapporti tra grandezze diverse.  
Esempi: altezza, peso, età, calorie.

*Remark 7.* I metodi per l'analisi delle variabili misurate su scale per intervallo o per rapporto *non differiscono* tra loro viceversa differiscono quelli da usati per quantitative discrete o qualitative.

## 7.3 Scrittura articolo

### 7.3.1 Autorship

ICMJE introduce quattro criteri per l'authorship di un lavoro [19]:

1. *contributo sostanziale alla concezione o disegno del lavoro; o alla acquisizione, analisi, o interpretazione dei dati;*

2. *scrittura o revisione* critica delle bozze per aspetti dal contenuto intellettuale rilevante;
3. *approvazione* finale della versione per la pubblicazione;
4. disponibilità ad essere *accountable* per tutti gli aspetti del lavoro, ad assicurare che domande relative ad accuratezza/integrità di qualsiasi parte del lavoro siano investigate e risolte.

La distinzione tra author e contributor raccomandata da ICMJE, avviene come segue:

**authors** coloro che rispettano tutti i 4 requisiti contemporaneamente;

**contributors** coloro ne rispettano solo alcuni, ma non tutti e 4<sup>1</sup>

Idealmente, come sempre ICMJE raccomanda, quando il lavoro di ricerca è condotto da un gruppo nutrito di persone, il gruppo stesso dovrebbe decidere *in anticipo* (prima che il lavoro inizi) chi figurerà come autore (incaricandolo del rispetto degli oneri dell'authorship) e confermare in sede di sottomissione chi lo debba essere veramente (alla luce del rispetto effettivo o meno degli oneri imposti dall'authorship di cui prima).

---

<sup>1</sup>Esempi di attività tipiche che, da sole, non sono sufficienti per l'authorship: supervisione di un gruppo di ricerca, supporto amministrativo, supporto di scrittura, editing tecnico, correzione bozze e lingua.



## Capitolo 8

# Confounding e interazione

```
## -----  
## WOODWARD standardizzazione tassi (esempio pag 178 sgg)  
## -----  
evs <- list('I' = c(0,0,1,6,7,16,17,25),  
            'II' = c(0,0,4,7,13,11,28,44),  
            'III' = c(0,0,1,9,17,19,43,53),  
            'IV' = c(0,1,5,10,15,24,28,56))  
popns <- list('I' = c(4784, 4210, 3396, 3226, 2391, 2156, 2182, 2054),  
              'II' = c(4972, 4045, 3094, 2655, 2343, 2394, 2597, 2667),  
              'III' = c(4351, 3232, 2438, 2241, 2360, 2708, 2968, 2802),  
              'IV' = c(4440, 3685, 2966, 2763, 2388, 2566, 2387, 2380))  
## standardizzazione diretta (non serve ev della popolazione standard)  
std_ev1 <- list('All' = rep(NA, 8))  
std_pop1 <- list('All' = c(8,6,6,6,6,5,4,4))  
  
## esempio per standardizzazione indiretta  
std_ev2 <- list('All' = c(0,1,11,32,52,70,116,178))  
std_pop2 <- list('All' = c(18547,15172,11894,10885,9482,9824,10134,9903))  
  
stdrate <- function(ev = NA,      # n. events per strata our sample  
                    pop = NA,     # n. pop (or exposure time) of our sample  
                    std_ev = NA,  # n. of events standard/reference pop  
                    std_pop = NA, # n. pop. (or exposure time) of std pop  
                    per = 1000,   # multiply rate per X  
                    ser_per = 100) # multiply ser per X  
{  
  ## ref woodward pag 181  
  if (length(pop) != length(std_pop))  
    stop("pop and std_pop must be of the same length.")  
  ## per il calcolo del crudo e standardizzazione diretta non serve std_ev  
  ## crude calculation  
  crude <- (sum(ev)/sum(pop)) * per  
  ## direct standardization  
  dstd <- sum((ev/pop)*std_pop ) / sum(std_pop)
```

```

dstd_se <-
  (1 / sum(std_pop)) *
  sqrt(sum(ev * (std_pop / pop)^2))
dstd_ci <- dstd + c(-1,1) * qnorm(0.975) * dstd_se
dstd_res <- setNames(c(dstd, dstd_se, dstd_ci),
  c('est', 'se', 'lower.ci', 'upper.ci')) * per

## ser/smr/sir ecc ecc ecc
exp_ev <- sum((std_ev / std_pop) * pop)
ser <- sum(ev) / exp_ev
ser_se <- sqrt(sum(ev)) / exp_ev
ser_res <- setNames(c(ser, ser_se), c('est', 'se')) * ser_per
## indirect standardization
crude_pop <- sum(std_ev) / sum(std_pop)
indstd <- ser * crude_pop
indstd_se <- crude_pop * ser_se
indstd_res <- setNames(c(indstd, indstd_se), c('est', 'se')) * per

## results
list('crude' = crude,
     'direct_std' = dstd_res,
     'ser' = ser_res,
     'indirect_std' = indstd_res)
}

Map(stdrate, evs, popns, std_ev1, std_pop1)

## $I
## $I$crude
## [1] 2.950941
##
## $I$direct_std
##      est      se lower.ci upper.ci
## 3.276607 0.388698 2.514773 4.038441
##
## $I$ser
## est se
## NA NA
##
## $I$indirect_std
## est se
## NA NA
##
##
## $II
## $II$crude
## [1] 4.320265
##
## $II$direct_std
##      est      se lower.ci upper.ci
## 4.1990966 0.4153992 3.3849292 5.0132641
##

```



```

## $II$ser
## est se
## NA NA
##
## $II$indirect_std
## est se
## NA NA
##
##
## $III
## $III$crude
## [1] 6.147186
##
## $III$direct_std
## est se lower.ci upper.ci
## 5.2993494 0.4615262 4.3947748 6.2039241
##
## $III$ser
## est se
## NA NA
##
## $III$indirect_std
## est se
## NA NA
##
##
## $IV
## $IV$crude
## [1] 5.896076
##
## $IV$direct_std
## est se lower.ci upper.ci
## 5.7544597 0.4933616 4.7874887 6.7214307
##
## $IV$ser
## est se
## NA NA
##
## $IV$indirect_std
## est se
## NA NA

(res1 <- Map(stdrate, evs, popns, std_ev2, std_pop2))

## $I
## $I$crude
## [1] 2.950941
##
## $I$direct_std
## est se lower.ci upper.ci

```

```

## 3.3795540 0.4001367 2.5953005 4.1638074
##
## $I$ser
##      est      se
## 69.718302 8.216381
##
## $I$indirect_std
##      est      se
## 3.3462108 0.3943547
##
##
## $II
## $II$crude
## [1] 4.320265
##
## $II$direct_std
##      est      se lower.ci upper.ci
## 4.3245158 0.4183908 3.5044850 5.1445466
##
## $II$ser
##      est      se
## 90.291430 8.728802
##
## $II$indirect_std
##      est      se
## 4.333642 0.418949
##
##
## $III
## $III$crude
## [1] 6.147186
##
## $III$direct_std
##      est      se lower.ci upper.ci
## 5.4252307 0.4576616 4.5282304 6.3222310
##
## $III$ser
##      est      se
## 113.028861 9.485171
##
## $III$indirect_std
##      est      se
## 5.4249513 0.4552518
##
##
## $IV
## $IV$crude
## [1] 5.896076
##

```

```

## $IV$direct_std
##      est      se lower.ci upper.ci
## 5.914942 0.502030 4.930982 6.898903
##
## $IV$ser
##      est      se
## 123.45627 10.47142
##
## $IV$indirect_std
##      est      se
## 5.9254268 0.5025881

## res2 <- Map(stdrate, evs, popns, std_ev2, std_pop2, list(1), list(1))
## res1[1]
## res2[1]

## -----
## mantel hanzel odd ratio, esempio pag 189
## -----

dimn <- list('housing' = c('rented', 'owner'),
             'chd' = c('yes', 'no'))
no_smoke <- matrix(c(33, 48, 923, 1722), ncol = 2, dimnames = dimn)
smoke <- matrix(c(52, 29, 898, 678), ncol = 2, dimnames = dimn)
library(lbmisc)
no_smoke_df <- table2df(as.table(no_smoke))
smoke_df <- table2df(as.table(smoke))
no_smoke_df$smoke <- 'No'
smoke_df$smoke <- 'Yes'

df <- rbind(smoke_df, no_smoke_df)
tab <- table(df)
mantelhaen.test(tab)

##
## Mantel-Haenszel chi-squared test with continuity correction
##
## data:  tab
## Mantel-Haenszel X-squared = 2.5354, df = 1, p-value = 0.1113
## alternative hypothesis: true common odds ratio is not equal to 1
## 95 percent confidence interval:
##  0.9537528 1.8203312
## sample estimates:
## common odds ratio
##      1.317629

## oppure, ma è forse più chiara la precedente ...
## with(df, mantelhaen.test(housing, chd, smoke))

```



Parte IV

Studi sperimentali



## Capitolo 9

# Studi sperimentali

**Definition 9.0.1** (Studio sperimentale). Studio in cui i ricercatori *assegnano* un trattamento (di varia natura)

**Example 9.0.1.** È passibile di studio sperimentale un farmaco o un intervento chirurgico; non lo è il fumo (somministrazione non etica) o il genere (somministrazione impossibile)

### 9.1 Fasi degli studi sperimentali (farmacologici)

**Definition 9.1.1** (Sperimentazione clinica (dlgs 1997)). Sperimentazione condotta *su soggetti umani* a verificare effetti clinici . . . di un prodotto in sperimentazione con l'obiettivo di valutarne sicurezza ed efficacia.

*Remark 8* (Fasi della sperimentazione). Lo sviluppo clinico di un farmaco è generalmente suddiviso in fasi temporali (dalla 1 alla 4), i cui confini non sono tracciabili in modo rigido [amadori2004sperimentazioneclinicaoncologia ]:

1. primi studi su un nuovo principio attivo condotti sull'uomo: scopo primario consiste nell'ottenere una valutazione preliminare su sicurezza e dosaggio del farmaco da impiegare (nonché avere le prime info su farmacocinetica e farmacodinamica nell'uomo)
2. scopo è dimostrare l'attività del farmaco (nonché valutarne meglio la sicurezza)
3. dimostrare l'efficacia terapeutica del farmaco
4. studi post commercializzazione per valutazione di sicurezza nel medio-lungo periodo

**Definition 9.1.2** (Attività di un farmaco). Capacità del trattamento di indurre le modificazioni della malattia grazie alle quali si presume che l'ammalato possa avere un beneficio.

**Definition 9.1.3** (Efficacia di un farmaco). Capacità del trattamento di indurre un beneficio clinico negli ammalati ai quali viene somministrato.





# Capitolo 10

## Fase 1

Costituiscono il primo passo nella sperimentazione sull'uomo dopo, il completamento degli studi preclinici (in vitro e successivamente in vivo, su animali).

### 10.1 Obiettivi

**Obiettivo primario** Determinare la dose di un farmaco da raccomandare per gli studi successivi, in relazione alla tossicità registrata.

*Remark 9.* Sebbene tossicità sia il termine comunemente utilizzato, “eventi avversi” è più appropriato dato che tossicità implica una relazione causale con il farmaco mentre “evento avverso” rimane più neutrale [13]

*Remark 10* (Criteri valutazione tossicità). Standard moderno è la classificazione CTCAE (Common Terminology for Adverse Events) che per ogni possibile evento avverso rilevato durante l'osservazione (ad esempio) definisce una scala di gravità crescente di severità da 1 (moderato) a 5 (decesso dovuto ad evento avverso).

**Example 10.1.1.** Per il CTCAE Term **Anemia**, il grado 1 corrisponde a Hemoglobin (Hgb)  $<LLN - 10.0$  g/dL;  $<LLN - 6.2$  mmol/L;  $<LLN - 100$  g/L.

**Obiettivi secondari** I principali sono descrivere farmacocinetica e farmacodinamica del trattamento.

**Definition 10.1.1** (Farmacocinetica (PK)). Modificazione nel tempo della concentrazione nel sangue di un farmaco.

*Remark 11.* È possibile costruire per una data dose e via di somministrazione una curva che rappresenti i valori di concentrazione plasmatica in funzione del tempo. Per tutte le vie di somministrazione ad eccezione dell'intravascolare l'andamento della curva prevede tre fasi: crescita (assorbimento farmaco), picco e decremento (metabolizzazione/smaltimento). Nel caso di somministrazione intravascolare non vi è assorbimento quindi il picco si ha all'inizio e la curva è monotona decrescente.

**Definition 10.1.2** (Farmacodinamica (PD)). Studio degli effetti biochimici e fisiologici del farmaco.

*Remark 12.* L'effetto si descrive costruendo un grafico che rappresenta l'entità della risposta (in termini ad esempio di espressione di una data proteina) in funzione del logaritmo della dose somministrata (perché si presuppone una relazione dose-risposta).

## 10.2 Popolazione

Sono generalmente condotti su volontari sani.

In oncologia, data l'utilizzazione di farmaci potenzialmente tossici, la sperimentazione non è fattibile su volontario sano (per la tossicità insita nel trattamento) ma è condotta su pazienti affetti da tumore che forniscono il proprio consenso.

*Remark 13.* Recentemente si è assistita alla crescita di interesse nei confronti dei cosiddetti phase 0 trial [13, p. 5] dove un *limitatissimo numero di soggetti sani* può essere arruolato per valutare microdosi di una nuova terapia al fine di rispondere a domande utili prima che una fase 1 vera e propria possa iniziare.

## 10.3 Definizioni

**Definition 10.3.1** (Dose massima tollerata (MTD)). Dose da raccomandare per gli studi delle fasi successive, coincidente con la massima dose per la quale la tossicità del farmaco (eventuali eventi avversi) risulta accettabile.

**Definition 10.3.2** (Dose massima somministrata (MAD)). Dose alla quale l'escalation cessa a causa dell'osservazione di un numero critico di DLT.

**Definition 10.3.3** (Tossicità dose-limitante (DLT)). Ogni evento tossico così severo o irreversibile tale da impedire l'incremento della dose.

*Remark 14.* La DLT è spesso definita come l'occorrenza di tossicità severa: gradi 3 o 4 per eventi avversi non ematologici o gradi 4 per tossicità ematologiche [13, p. 171].

**Definition 10.3.4** (Recommended phase II dose (RP2D)). Dose raccomandata per la fase 2 (concide con la MTD)

## 10.4 Disegni

Uno studio di fase 1 è tipicamente disegnato come uno studio a dosi crescenti per la determinazione della MTD. In merito alla dose somministrata le decisioni da prendere sono:

- la dose (espressa in milligrammi per metro quadrato di superficie corporea) dalla quale partire: scelta comune consiste nel 10% della MELD<sub>10</sub> (o  $LD_{10}$ )

**Definition 10.4.1** (MELD<sub>10</sub> (o  $LD_{10}$ )). Dose letale per il 10% dei topi sottoposti a tale dose.

- lo schema di incremento: trade off tra aumento troppo veloce (esporre pazienti a tossicità eccessive) o troppo lento (allungamento dei tempi di sviluppo di un farmaco potenzialmente utile). Ogni disegno (vedi in seguito) propone un diverso schema.

L'aumento della dose è possibile solo dopo che è trascorso un periodo di tempo sufficientemente prolungato per osservare l'eventuale effetto tossico nei pazienti inseriti a livello precedente.

I disegni di fase 1 possono essere classificati, in base allo schema di incremento, in due gruppi [13]:

**rule-based o algoritmici** determinano la dose attraverso un processo iterativo: i pazienti vengono assegnati a dose via via crescenti in base a regole prespecificate in protocollo, in relazione ad eventi di tossicità.

**Example 10.4.1.** Il disegno standard 3+3

**model-based** stima la relazione dose tossicità e procede alla scelta della dose

**Example 10.4.2.** Continual reassessment method

### 10.4.1 Disegno standard (3+3)

#### 10.4.1.1 Funzionamento

Il disegno standard utilizza:

- un profilo di incrementi decrescenti del livello di dose secondo lo schema di tabella 10.1; in particolare la dose non viene mai variata nel singolo paziente ma si considerano nuovi pazienti ai quali somministrare la nuova dose;
- un algoritmo di passaggio a dosaggi via via successivi come da diagramma 10.1.

La procedura, dal punto di vista statistico si basa sulle seguenti considerazioni:

- se almeno 2 su 3 pazienti trattati ad un particolare livello di dose mostrano una DLT, si può affermare con confidenza del 90% che la probabilità di DLT a quella dose è  $> 20\%$

```
binom.test(2, 3, alternative = 'greater', conf.level = 0.9)$conf.int
## [1] 0.1958001 1.0000000
## attr(,"conf.level")
## [1] 0.9
```

- d'altra parte, se 0 pazienti mostrano una DLT si può affermare con confidenza del 90%, che la vera probabilità di DLT è  $< 55\%$

```
binom.test(0, 3, alternative = 'less', conf.level = 0.9)$conf.int
## [1] 0.0000000 0.5358411
## attr(,"conf.level")
## [1] 0.9
```

Livello	Aumento	Esempio dose ( $\text{mg}/\text{m}^2$ )
1	10% del MELD10	1
2	100	2
3	67	3.3
4	50	5.0
5	40	6.7
6	33	8.8
7	33	11.8
8	33	15.7

Tabella 10.1: Schema incrementi di dose. Da [amadori2004sperimentazioneclinicaoncologia].

Nel caso lo studio non riesca a identificare la MTD si può procedere in due direzioni:

- proporre la dose più elevata come dose ottimale per la fase 2
- progettare una nuova fase 1 che preveda dosi ancor più elevate di farmaco

Nel caso si esca per MTD identificata alla dose iniziale occorre disegnare un nuovo studio di fase 1 che parta da una dose iniziale inferiore.

#### 10.4.1.2 Critiche

Le principali [amadori2004sperimentazioneclinicaoncologia]:

- troppi pazienti vengono trattati a basse dosi, con poca utilità da un punto di vista terapeutico
- incrementi lenti e numerosi, durata lunga per la definizione della MTD

#### 10.4.1.3 Varianti dello schema

**Dose iniziale maggiore** Es impostando la dose iniziale al 20% della MELD invece che al 10%

**Accelerated titration design** L'idea è trattare un primo paziente a livelli di dose via via crescenti fino a quando non si osserva una tossicità di grado 2 o maggiore. Quando questo avviene si trattano i successivi pazienti partendo dalla dose precedente e seguendo lo schema classico.

### 10.4.2 Continual reassessment method (Adattamento continuo)

Si stima in maniera Bayesiana la relazione tra dose e tossicità per derivare il livello di dose che è associato ad una determinata frequenza di DLT (solitamente il 20-30%):

- si parte da una distribuzione a priori della probabilità di DLT in funzione della dose

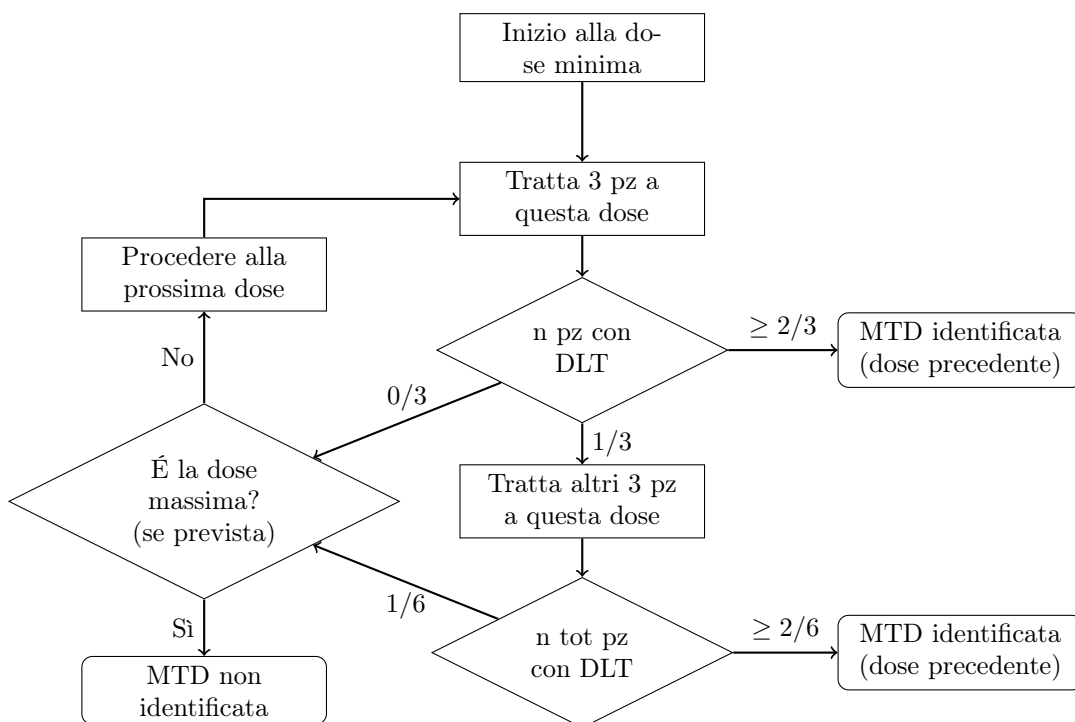


Figura 10.1: Flowchart schema classico 3+3. Da [13].

- ogni nuovo paziente che entra nello studio viene trattato alla dose stimata come MTD (allo stato attuale)
- la risposta del paziente viene utilizzata per aggiornare il modello in maniera tale da aggiornare la distribuzione a posteriori della relazione dose/tossicità sulla base della quale effettuare la scelta di dose per il pz successivo

Con questo metodo:

- occorre inserire un criterio di uscita dallo studio ma inserendo nello studio 20-25 pazienti si ottiene una stima sufficientemente accurata di MTD
- ogni paziente viene trattato con una dose di farmaco sufficientemente alta (con un beneficio clinico per tutti non solo gli ultimi)

Per ulteriori approfondimenti: [13, p. 178] e [24].

### 10.4.3 Disegno per nuovi farmaci

Laddove non vi sia correlazione tra dose somministrata ed effetto cade il paradigma sul quale si fonda il disegno classico di fase 1 (utile per lo più per farmaci tipo i citotossici).

L'idea è che la dose ottimale da mandare in fase 2 è una non necessariamente tossica e procedere in maniera classica porterebbe ad un sovradosaggio del farmaco.

**TODO:** Verifica disegni per nuovi farmaci



# Capitolo 11

## Fase 2

Screenano farmaci poco attivi avviando a nuovi studi di più ampie dimensioni soltanto quelli potenzialmente efficaci.

### 11.1 Obiettivi

**Obiettivo primario** Valutare l'attività del trattamento intesa come capacità di indurre le modificazioni sulla malattia/paziente che fanno ben sperare per la sua prognosi.

**Obiettivi secondari** Principalmente tossicità, ulteriori analisi di farmacocinetica e farmacodinamica, valutazione effetti su specifici bersagli molecolari (se noti)

### 11.2 Aspetti da considerare nel disegno

Alcuni aspetti da chiarire nel disegno di uno studio di fase 2 ([7] e [amadori2004sperimentazioneclinicaoncologi] seguono.

#### 11.2.1 Popolazione

Uno studio di fase 2 richiede pazienti:

- affetti da uno specifico tipo di neoplasia
- per i quali si abbia avuto progressione dopo una terapia standard
- per i quali sia possibile valutare la risposta della malattia al trattamento (es dimensione tumore o altri indicatori)
- per i quali il trattamento possa dare un beneficio

Per questo motivo:

- nel caso di farmaci citotossici sono esclusi che presentano una malattia non misurabile (poiché non è possibile definire una risposta parziale);

- possono essere esclusi pz con performance status basso e/o aspettativa di vita  $< 3$  mesi;
- i pz non debbono avere gravi malattie concomitanti
- adeguata funzionalità renale, epatica, cardiaca e respiratoria.
- possono essere esclusi pz con precedenti trattamenti (non standard, altri studi): potrebbero debilitare il paziente rendendo difficile poter somministrare la dose piena del farmaco sperimentato;

### 11.2.2 Trattamento

Quale è il *meccanismo d'azione*? Questa avrà influsso su diversi aspetti, principalmente l'outcome adottato e l'impiego della randomizzazione o meno.

Il trattamento è (o è assimilabile a) un:

- **citotossico**: l'obiettivo del trattamento (o comunque come è ragionevole misurare il funzionamento del farmaco) è la contrazione della massa tumorale.

Se la risposta al trattamento è misurabile attraverso la dimensione del tumore i criteri standard attuali sono i RECIST [12].

Data l'estrema rarità di una regressione tumorale spontanea, è giustificato ipotizzare che la risposta che si è eventualmente verificata possa essere solo la conseguenza del farmaco e ciò rende possibile la programmazione di studi a singolo braccio.

- **citostatico**: l'obiettivo del trattamento non è tanto una diminuzione del volume, piuttosto evitare/rallentare la progressione.  
Nel caso di farmaci non citotossici il criterio dimensionale su cui si basa la valutazione della risposta dei trattamenti classici viene tipicamente a cadere. Si rende necessaria l'identificazione di nuovi parametri per valutare l'attività del trattamento [amadori2004sperimentazioneclinicaoncologia ]:

- parametri clinici: valutazioni temporali quali tempo alla progressione, sopravvivenza, numero di pazienti vivi o liberi da progressione
- criteri biologici: ad esempio parametri misurabili su campioni biotici o su altro materiale biologico

Diviene necessario uno studio controllato necessario laddove si utilizzano indicatori che possono modificarsi spontaneamente nel tempo.

### 11.2.3 Outcome

L'outcome deve essere un surrogato validato, ossia tipicamente correlato ad un outcome hard/di efficacia usato in fase 3. Outcome tipici sono:

- risposta dicotomica (es CR+PR e parziale vs SD+PD)
- risposta multinomiale (es CR+PR vs SD vs PD)
- quantitativa (es marker)
- tempi all'evento:



- TTP: dalla randomizzazione alla progressione
- TTF: dalla randomizzazione all'interruzione del trattamento (per progressione, tossicità, scelta del paziente o morte)
- PFS: dalla randomizzazione alla progressione o morte (quale
- tempo ad un evento rilevante (es al SAE per studi concentrati sulla tossicità, alla prima frattura per trattamenti contro metastasi ossee)

#### 11.2.4 Randomizzazione

La randomizzazione è necessaria se:

- l'outcome scelto è soggetto ad evoluzione spontanea
- dati storici non sono disponibili o sono poco affidabili/confrontabili eccetera
- molteplici trattamenti oppure molteplici dose/schedule/sequenze di un singolo trattamento

#### 11.2.5 Scopo (sottofase) dello studio

A volte si suddividono gli studi di fase II in due sottogruppi:

1. phase IIa (*proof of concept*): si valuta l'attività di un trattamento che ha completato la fase 1 oppure si valutano diverse dosi di un trattamento per valutare la dose-risposta (learning trial). Più tipicamente a singolo braccio.
2. phase IIb (*go/no-go*): studio nel quale si decide se procedere o meno con la fase 3, può comparare diversi dosaggi o un dosaggio verso placebo, tipicamente randomizzando.

A volte le due sottofasi vengono riunite in un unico studio.

#### 11.2.6 Categorie disegni

**Stadio unico** Un campione prefissato di pazienti è reclutato, trattato e seguito sino all'ottenimento dell'outcome. Questi disegni:

- sono semplici: evitano complessità inerenti strategie di reclutamento in presenza di analisi ad interim
- non permettono flessibilità, es fermare il trial in anticipo a causa di livelli troppo bassi (o veramente alti) di attività

Sono utili se:

- il trattamento è molto efficace e si presume di concludere lo studio con pochi pazienti;
- vi è relativa sicurezza sulla tossicità tale da non rendere necessaria una analisi ad interim.

**Due stadi** Sono opportuni se l'attività del trattamento è sconosciuta e/o la tossicità non può essere sottovalutata (ci si lascia spazio per terminare in anticipo).

I pazienti sono reclutati in due fasi:

- al termine della prima si decide se continuare o fermarsi (mancanza di attività o forte presenza di attività, tossicità inaccettabile) oppure su quale trattamento portare alla seconda fase;
- al termine della seconda si fornisce la stima di attività

Data la natura del disegno il sample size complessivo non è fisso, pertanto di solito questi disegni forniscono un sample size massimo e un sample size medio (ASN, che contempra la probabilità di terminazione precoce e i sample size nei due stati del mondo).

**Multistage** Anche detti group-sequential, sono simili ai disegni a due stadi (ma ad un numero di stadi maggiori) e vengono specificati criteri di stop per ciascuno di essi.

**Continuous monitoring** La valutazione di attività è effettuata ogni volta che un nuovo paziente ha il dato di attività (anche qui lo si fa per early termination alla massima potenza).

**Decision theoretic** Valutano costi e guadagni associati ad effettuare decisioni sbagliate alla fine di una fase 2, incorporando funzioni di utilità associate a questi costi/guadagni. Valutazione più generale rispetto a quella meramente clinica, ma di difficile applicazione; raramente usati in oncologia.

**Three outcomes** Alla fine dello studio non vi sono solo due ipotesi (nulla e alternativa) ma si può ottenere anche uno studio non conclusivo

**Phase 2/3** Sono usati quando una transizione dalla fase 2 alla 3 deve essere il più efficiente possibile: permettono di incorporare i dati generati dalla fase 2 nell'analisi di fase 3 (risparmiando così pazienti). Sono solitamente randomizzati con gruppo di controllo. La randomizzazione può essere usata per:

- decidere se continuare un singolo trattamento da portare in fase 3;
- selezionare il trattamento con la miglior risposta dei diversi confrontati, da portare in fase 3.

Per evitare problemi di molteplicità, nell'analisi inerente la fase 2 si usa un outcome diverso da quello di fase 3 (e si inizia a misurare l'outcome di fase 3 anche nei pazienti che iniziano con la 2; chiaramente facendo i due studi in maniera separata si può beneficiare della conoscenza accumulata durante la fase 2 nel disegno della fase 3 (es cambi ai criteri di eleggibilità, schema follow up).

**Randomised discontinuation** Nei trial a randomized discontinuation (o *enrichment trial*) i pazienti vengono trattati tutti con lo stesso trattamento e in seguito valutati per risposta

- quelli in progressione escono dallo studio
- quelli in risposta (parziale/completa) continuano a ricevere il trattamento
- quelli in stable disease sono randomizzati a continuare il trattamento, a interromperlo (e rimanere senza trattamento) oppure a ricevere un trattamento standard (in relazione alla domanda dello studio)

## 11.3 Disegni comuni

Nell'approccio frequentista si ha un test di ipotesi unilaterale del tipo

$$\begin{aligned} H_0 : p &\leq p_0 \\ H_1 : p &\geq p_1 \end{aligned}$$

con:

- $p_1 > p_0$
- $p_0$  livello di risposta giudicato fallimentare (oppure quella del trattamento standard se disponibile)
- $p_1$  un livello di risposta auspicabile;  $p_1$  è scelto in modo tale che  $p_1 - p_0$  rappresenti il miglioramento auspicabile col nuovo trattamento

### 11.3.1 Livelli di errore nell'inferenza

Ponendoci in ottica di screening di potenziali candidati, riguardo agli errori:

- l'errore del primo tipo,  $\alpha$ , può andare dal 5% tipico della fase 3 sino al 20%, al fine di risparmiare pazienti;
- l'errore di tipo 2,  $\beta$ , deve comunque rimanere al massimo 20% (meglio 10%); si vuole comunque evitare di avere dei falsi negativi in questa fase di screening, similmente a quanto avviene nella fase 3.

### 11.3.2 Stadio unico - A'Hern

Adottando una distribuzione binomiale dei successi si giudica la combinazione ampiezza campionaria/eventi osservati soddisfacente

$$\mathbb{P}(\text{reject } H_0 | \pi = p_0) \leq \alpha \mathbb{P}(\text{accept } H_0 | \pi = p_1) \leq \beta$$

A livello di algoritmo

- si scelgono i parametri  $\alpha, \beta$ ;
- si cicla sul sample size (diciamo da 1 a 2000);
- per ogni sample size si cicla sul numero di successi (da 0 a sample size) a 2000), considerati come soglia per rifiutare la nulla;

- si calcolando le probabilità associate alle binomiali sotto nulla e alternativa;
- si ritiene la combinazione sample size/eventi con sample size minore

**Example 11.3.1.** `lbss::ahern(p0 = 0.4, p1 = 0.55, alpha = 0.1, beta = 0.1)`  
`## Error in loadNamespace(x): non c'è alcun pacchetto chiamato 'lbss'`

Per impostare il trial come sopra che giudichi il farmaco utile per la fase 3 se consente un incremento di 15 punti percentuali nella risposta si ha che

- occorre reclutare 75 pazienti
- si può passare alla fase 3 se almeno 36 vanno in risposta; infatti in tal caso il limite inferiore dell'intervallo di confidenza (Clopper Pearson, al 90% ad una coda), è 0.4006 che è maggiore di  $p_0$ .

### 11.3.3 Due stadi - Simon

È un disegno:

- a due stadi
- per outcome binario
- che permette early termination per mancanza di attività (non è previsto early stop per eccesso di attività)

Se:

- $n_1$  sono i pazienti da reclutare al primo stage
- $r_1$  i successi che fungono da cutoff al primo stage ( $r_1$  si termina  $r_1 + 1$  si prosegue)
- $n_2$  i pazienti da reclutare al secondo stage
- $r$  i successi che fungono da cutoff al secondo stage (se  $r$  lo studio non va in fase 3 se  $r + 1$  sì)

Le seguenti sono la probabilità di early (al primo step) termination PET, il numero atteso di pazienti EP, la probabilità di rifiutare  $H_0$  al termine dello studio (R) (serve per calcolare la potenza direi se)

$$PET(p) = B(r_1; p, n_1)$$

$$EP(p) = n_1 + (1 - PET(p)) \cdot n_2$$

$$R(p) = 1 - B(r_1; p, n_1) - \sum_{m=r_1+1}^{\min(n_1, r)} b(m; p, n_1) \cdot B(r - m; p, n_2)$$

dove  $b$  è la PMF e  $B$  la cumulative distribution function della binomiale.

Le quantità di sopra possono essere calcolate sotto varie ipotesi di  $p$  ma di interesse nell'ipotesi che sia vera  $H_0$  ( $p = p_0$ ).

Del disegno di [22], una volta fissati  $\alpha, \beta$  e indagando per forza bruta ve ne sono due varianti:

- **optimal**: minimizza i pazienti attesi sotto  $H_0$  ( $EN(p_0)$ ) nonchè i pazienti da reclutare al primo stage
- **minimax** minimizza i pazienti complessivi ( $n_1 + n_2$ ) al termine delle due fasi (da preferire nel caso vi possano essere difficoltà di reclutamento, es tumori rari)

**Example 11.3.2.** Se  $p_0 = 0.4$ ,  $p_1 = 0.55$ ,  $\alpha = \beta = 0.1$  e si considererà interessante il farmaco qualora permetta un incremento di .15 rispetto alla risposta standard.

```
# pu è p0, pa è p1, ep1 è alpha, ep2 è beta
clinfun::ph2simon(pu = 0.4, pa = 0.55, ep1 = 0.1, ep2 = 0.1)

## Error in loadNamespace(x): non c'è alcun pacchetto chiamato 'clinfun'
```

Lo studio si svolgerà come segue:

- durante la prima fase bisogna arruolare 38 pazienti con Optimal (45 con Minimax)
- si termina alla prima fase (farmaco *non* di interesse) se le risposte sono al massimo 16 (18 per minimax), altrimenti si prosegue alla seconda fase
- durante la seconda fase bisogna arruolare altri 50 pazienti (28 per minimax), per arrivare complessivamente a 88 pazienti (73 per minimax)
- il farmaco non è di interesse se le risposte sono complessivamente (fase 1 e 2) al massimo 40 (o 34 per minimax), altrimenti si può proseguire con studi di fase 3
- PET: probabilità di terminare lo studio allo step iniziale sotto  $H_0$  ( $p = p_0$ ) è 66% (56% per minimax) ed è calcolata mediante la cumulata della binomiale di  $r_1$  risposte su  $n_1$  trial sotto ipotesi che  $p = p_0$  (terminiamo al primo stage se osserviamo  $r_1$  risposte o meno), ossia

```
pbinom(16, 38, 0.4)

## [1] 0.6695864
```

- EN: il numero atteso di pazienti da reclutare sotto l'ipotesi  $H_0$  è 54 (57 per minimax) ed è calcolato come

$$EN = n_1 + (1 - PET)n_2$$

con  $n_1$  pazienti della fase iniziale,  $n_2$  della seconda fase e PET probabilità di terminare alla fine della prima fase.

```
38 + (1 - pbinom(16, 38, 0.4)) * 50

## [1] 54.52068
```

**Example 11.3.3.** Provando a riprodurre a mano l'algoritmo: fissato  $\alpha, \beta, p_0, p_1$

- for sample size tot:  $n \in (2, \dots, max\_samplesize)$
- for sample size tot:  $n_1 \in (1, \dots, n - 1)$
- ottieni sample size fase2:  $n_2 = n - n_1$
- for successi possibili fase1  $r_1 \in (0, n_1)$
- for successi possibili fase 2: *what*
- calcola la probabilità di rifiutare l'ipotesi nulla sotto ipotesi che sia vera (e verifica che tale valore sia inferiore a  $\alpha$ ) o sia vera  $H_1$  (e che sia maggiore di  $1 - \beta$ )
- restituisci le combinazioni di parametri che rispettano  $\alpha, \beta$

### 11.3.4 Altri disegni

Si hanno:

- [15] facilita l'early stopping rispetto a Simon per mancata attività, mentre [10] rimuove alcune limitazioni
- [9] permettono l'early stopping sia in caso di bassa che alta attività
- [8] considerano congiuntamente l'analisi di attività e tossicità, i due obiettivi principali di uno studio di fase II

Per un confronto e un aiuto nella scelta, complessivamente, [20].

### 11.3.5 Stima al termine di un multistage

Occorre adottare procedure adhoc come spiegato in [20]

## 11.4 Criteri RECIST

### 11.4.1 Classificazione delle lesioni e tumour burden

Al baseline lesioni/linfonodi sono classificati come:

- *misurabile*:
  - lesioni tumorali con diametro maggiore  $\geq 20\text{mm}$  (se misurato con raggi) o  $\geq 10\text{mm}$  (CT scan, calibro)
  - linfonodi con asse minore  $\geq 15\text{mm}$  (CT scan)
- *non misurabile*: tutte le rimanenti lesioni

Il metodo di misurazione impiegato al baseline deve rimanere costante ai successivi follow-up.

**Definition 11.4.1** (Measurable disease). Presenza di almeno una lesione misurabile.

*Remark 15.* In protocolli [12, p. 232]:

- in cui la risposta è l'endpoint primario solo i pazienti con measurable disease al baseline dovrebbero essere inclusi in protocollo;
- in cui la progressione tumorale è l'outcome primario (es pfs o proporzione di pazienti con progressione ad un dato istante) il protocollo deve specificare se l'ingresso in analisi è ristretto a coloro che hanno measurable disease o tutti (anche a pazienti con nemmeno una lesione misurabile).

*Remark 16.* Al fine di determinare l'*overall tumor burden* al baseline, le lesioni vengono ulteriormente suddivise in lesioni target e non

**Definition 11.4.2** (lesioni target). lesioni misurabili, fino ad un massimo di 2 per organo e di 5 in totale nel singolo paziente, scelte in modo da risultare rappresentative di tutti gli organi interessati; debbono essere scelte in base alla dimensione del diametro maggiore e all'accessibilità prevedibile nelle successive valutazioni;

**Definition 11.4.3** (lesioni non target). tutte le lesioni rimanenti.

**Definition 11.4.4** (Dimensione al baseline (overall tumor burden)). Somma dei diametri maggiori (asse minore per i linfonodi) di tutte le lesioni target, rilevata all'inizio del trattamento

*Remark 17.* La dimensione al baseline **costituisce il riferimento** per le successive misurazioni e per la valutazione della risposta. La misura deve essere monitorata lungo il follow up per la determinazione della risposta

## 11.4.2 Risposta

*Remark 18.* Per valutare la risposta globale del paziente si valuta:

- risposta nelle lesioni target
- risposta nelle rimanenti
- comparsa di nuove lesioni

**Definition 11.4.5** (Risposte lesioni target). Sono:

- *risposta completa (CR)*: scomparsa di tutte le lesioni target, tutti i linfonodi patologici (target o non target) con asse inferiore  $< 10\text{mm}$ ;
- *risposta parziale (PR)*: diminuzione  $\geq 30\%$  rispetto alla dimensione al baseline della somma dei diametri delle lesioni target
- *progressione (PD)*: aumento  $\geq 20\%$  (ma comunque  $\geq 5\text{mm}$ ) rispetto rispetto alla dimensione minima registrata durante il follow-up (sia essa quella del baseline o meno);
- *stabilità di malattia (SD)*: ne riduzione in grado di dare PR, ne incremento tale da qualificare PD (caso residuale rispetto ai precedenti due).
- nel caso non tutte le lesioni target abbiano misurazione ad un dato follow up la risposta si considera *non valutabile (NE)*.

Lesioni target	Lesioni non target	Nuove lesioni?	Overall response
CR	CR	No	CR
CR	Non-CR/non-PD	No	PR
CR	Non valutate	No	PR
PR	Non-PD o non tutte valutate	No	PR
SD	Non-PD o non tutte valutate	No	SD
Non tutte valutate	Non-PD	No	NE
PD	Qualsiasi	Sì o No	PD
Qualsiasi	PD	Sì o No	PD
Qualsiasi	Qualsiasi	Sì	PD

Tabella 11.1: Valutazione risposta globale (OR) nei pz con malattia misurabile

**Definition 11.4.6** (Risposte lesioni non target). Sono:

- *risposta completa (CR)*: scomparsa di tutte le lesioni non target e negativizzazione dei marcatori tumorali. Tutti i linfonodi non patologici (asse minore < 10 mm);
- *Non-CR/Non-PD*: persistenza di una o più lesioni non target e/o persistenza di livelli elevati dei marcatori tumorali;
- *progressione (PD)*: inequivocabile progressione di lesioni non target preesistenti e/o comparsa di 1 o più lesioni

**Definition 11.4.7** (Risposta globale (overall response, OR)). Risposta complessiva del paziente al trattamento, che contempla sia lesioni target che non; nel caso di pazienti con malattia misurabile al baseline è definita sulla base di tabella 11.1. Per altri casi cfr [12, p. 235].

### 11.4.3 Outcome derivabili

**Definition 11.4.8** (Migliore risposta globale/complessiva (best overall response)). È la migliore risposta registrata dall'inizio sino alla fine del trattamento

**Example 11.4.1.** Un paziente che abbia SD alla prima valutazione, PR alla seconda e PD all'ultima ha come miglior risposta globale PR

*Remark 19.* La miglior risposta complessiva in genere è l'indicatore più utilizzato per definire l'attività di un farmaco in una fase 2.

**Definition 11.4.9** (Durata della risposta complessiva). Misurata dal momento in cui si ha CR o PR (quale che sia registrata prima) fino alla prima data in cui si sia documentata una PD o una ripresa di malattia

**Definition 11.4.10** (Durata della risposta completa). Dal momento in cui per la prima volta si ha CR fino alla prima data di recurrent disease

**Definition 11.4.11** (Durata della stabilità di malattia). Misurata dall'inizio del trattamento sino a che è registrata una PD



## Capitolo 12

# Feasibility/Pilot studies

### 12.1 Definizioni

In letteratura ci sono almeno due accezioni, soprattutto per fattibilità

- gli studi pilota/di fattibilità sono trial randomizzati piccoli che vengo fatti per vagliare la fattibilità di un trial futuro propriamente dimensionato. La *fattibilità* studiata è quella *dello studio*. L'estensione del consort [14] va ovviamente in questa direzione
- nell'ambito della ricerca sulla implementazione di un trattamento si può studiare la fattibilità dello stesso. La *fattibilità* studiata è quella *del trattamento*

### 12.2 Approccio a Maglietta - (mail mia 12/1/23)

Cari Debora e Silvio, nell'intento di condividere qualcosa di eventualmente utile ai fini dei rapporti con lo statistico del CE (Maglietta, AO PR), vi riporto come sto impostando ultimamente studi (approvati e) dimensionati in base a fattibilità (sicuramente i retrospettivi, forse anche i prospettici, se es si vuole andare avanti al max tot tempo a reclutare).

Per quella che è la mia esperienza non è che al CE guardino necessariamente male/non accettino gli studi dimensionati con quel che c'è/riusciamo a fare; solo, affinché non siano contestati su sto punto, bisogna prestare particolare attenzione al wording degli obiettivi dello studio (a parità di analisi statistica eseguita).

In sintesi, il canovaccio che personalmente adotterò laddove figuri come responsabile statistico (fino a "fail"/prova contraria), è circa il seguente:

- Sample size: (solita frase) fatto secondo fattibilità di rilevazione/in assenza di ipotesi a priori da sottoporre a verifica, variata in base alle circostanze di studio/ispirazione giornaliera
- Obiettivo primario: usare i termini "descrizione"/"esplorazione", NO "analisi"/"valutazione"

- Analisi primaria: analisi inferenziale con magari maggiore enfasi su intervalli di confidenza più che su test/verifica di ipotesi

Cose che mi fanno pensare che questo possa essere un approccio utile:

- la mia esperienza a Parma (Maglietta è stato quello che mi ha sostituito (dopo un po' di prove con altra gente) quando sono venuto a Reggio) e la filosofia ivi dominante: a volte con un feticcio eccessivo per il formalismo/"wording", ben più di quanto un anglosassone (il referee, nostro target di riferimento) abbia in mente;
- due studi osservazionali approvati utilizzando questo approccio (caso1.pdf protocollo finale; caso2.pdf richieste modifiche e protocollo con track changes che è stato in seguito approvato) e uno studio sperimentale prospettico (caso3.pdf, cfr primo/terzo punto) in attesa di valutazione in cui la richiesta di modifica del CE è in linea. Ovviamente, non c'è bisogno di dirlo, ma vi allego il tutto per vs eventuale consultazione/utilizzo no condivisione bla bla bla..

Se può per caso esserVi eventualmente utile mi fa piacere.  
un saluto, Luca

# Capitolo 13

## Fase 3

### 13.1 Obiettivi

Obiettivo *primario* è valutare l'efficacia del trattamento (nel prolungare la sopravvivenza rispetto allo standard attuale); obiettivi *secondari* sono:

- valutare ulteriormente la *risposta*
- valutare ulteriormente la *tollerabilità/eventi avversi*
- impatto sulla *QOL*

### 13.2 Classificazione di studi

#### 13.2.1 Studi esplicativi e pragmatici

*Remark 20.* Questa classificazione presenta profondo legame con il concetto di validità esterna e dirette implicazioni nell'analisi statistica (ITT vs PP).

**Definition 13.2.1** (Studi esplicativi (Explanatory trials)). Volti alla dimostrazione dell'efficacia del trattamento in *condizioni ideali*

**Definition 13.2.2** (Studi pragmatici). Mirati a valutare l'efficacia in un contesto assistenziale reale

### 13.3 Validità di uno studio

*Remark 21.* Uno studio è tanto più utile quanto più valido. Ve ne sono due accezioni

**Definition 13.3.1** (Validità interna). Grado con cui i risultati di uno studio sono vorretti per i pazienti che ne fanno parte (assenza di bias, dipende dalla conduzione)

**Definition 13.3.2** (Validità esterna). Generalizzabilità delle stime ad altri contesti

### 13.3.1 Validità interna

*Remark 22.* Uno studio è “valido internamente” se nella pianificazione e conduzione non si sono verificati tre tipi di bias (Sackett):

- Distorsione da selezione (selection bias): sbilanciamento tra i trattamenti nella distribuzione di fattori capaci di influenzare l’end-point.
- Distorsione di valutazione (assessment bias): sbilanciamento tra i trattamenti nel modo in cui i soggetti sono seguiti/valutati nel corso dello studio
- Distorsione di analisi (analysis bias): distorsione che avviene in fase di analisi dei dati

*Remark 23.* La validità interna è il punto di forza dei disegni sperimentali (soprattutto RCT) rispetto agli osservazionali.

### 13.3.2 Validità esterna

*Remark 24.* Uno studio è “valido esternamente” se i suoi risultati (siano essi distorti o meno) sono generalizzabili ad altri contesti.

*Remark 25.* E’ la domanda che implicitamente si pone un clinico quando valuta un RCT altrui per scelte terapeutiche nei confronti di propri pazienti.

*Remark 26.* La generalizzabilità è legata a:

- criteri di inclusione/esclusione dei pazienti
- setting dove lo studio è condotto (Es Ospedali Hi-Tech vs. arruolamento “sul territorio”)
- principi di analisi impiegati nella stima (ITT, PP)

## 13.4 PICO

### 13.4.1 Popolazione

#### 13.4.1.1 Criteri di inclusione/esclusione

**Definition 13.4.1** (Principio di incertezza (*equipoise*)). La scelta della popolazione deve essere fondato su questo: definire i criteri di inclusione/esclusione (su paziente/malattia) individuando coloro per i quali il medico è indeciso su quale possa essere il miglior trattamento.

*Remark 27.* La popolazione è definita in relazione agli obiettivi: un approccio pragmatico può portare ad allargare le maglie, mentre uno esplicativo a contrarle.

### 13.4.1.2 Popolazione d'analisi

Per quanto accuratamente condotta la ricerca, è quasi inevitabile un qualche scostamento dal protocollo per alcuni pazienti; ad esempio:

- pazienti inseriti che non rispettano i criteri di eleggibilità
- pazienti nel gruppo di trattamento che non hanno completato lo stesso, o pazienti del gruppo di controllo che hanno ricevuto trattamento sperimentale

In sede di analisi *di efficacia* si pone se includere o meno le info dei pazienti non aderenti al protocollo, e questo dipende dalla tipologia di approccio adottato nello studio (pragmatico o esplicativo)

**Approccio pragmatico e popolazione ITT** Se l'obiettivo è accertare il beneficio del trattamento in reali condizioni di pratica clinica, le deviazioni dal protocollo ricreano all'interno dello studio condizioni vicine alla pratica clinica; pertanto secondo questo approccio l'analisi deve essere condotta in base al principio dell'intenzione a trattare (ITT); l'insieme di tutti i soggetti randomizzati, aderenti o meno al protocollo, forma la cosiddetta popolazione intention-to-treat e su questa si concentrano le analisi.

**Approccio esplicativo e popolazione PP** Se si vuole valutare l'efficacia nelle condizioni ideali del trattamento l'analisi si deve limitare a quei pazienti che sono aderenti al protocollo, che nel loro insieme formano la popolazione per protocol. Vanno stabiliti i criteri minimi sufficienti per giudicare il trattamento ricevuto come adeguato.

**Popolazione di safety** Per quanto riguarda la valutazione *di tollerabilità* del trattamento ci si basa sull'insieme dei pazienti che hanno ricevuto almeno una dose, indipendentemente da qualunque altro fattore, e questa è la *safety population*.

## 13.4.2 Outcome

Gli outcome più utilizzati (soprattutto in ambito oncologico) sono:

- Overall Survival
- Progression Free Survival
- Time to progression

## 13.5 Randomizzazione

### 13.5.1 Alcuni concetti

**Definition 13.5.1** (Blinding). Consiste nella non conoscenza del braccio di allocazione da parte di un determinato gruppo di soggetti

*Remark 28.* Tradizionalmente gli studi possono essere classificati come:

- in *aperto*: nessun soggetto è cieco rispetto all'allocazione
- in *singolo cieco*: il paziente non è a conoscenza del braccio
- *doppio cieco*: sia paziente che medico che somministra/effettua il trattamento non sono a conoscenza del braccio
- *triplo cieco*: paziente, medico che tratta (es oncologo) e valutatore di outcome se differente (es radiologo) non sono a conoscenza dell'outcome

Spesso meglio esplicitare per esteso i soggetti che sottostanno a blinding

*Remark 29.* Serve per eliminare possibili bias nella valutazione d'outcome:

- il paziente potrebbe avere un effetto benefico non dovuto al trattamento nel sapere di essere stato trattato nel braccio sperimentale;
- chi somministra il trattamento (se non all'oscuro del braccio) potrebbe suggerirlo involontariamente al paziente;
- chi valuta potrebbe inconsciamente o meno biasare le valutazioni.

*Remark 30.* Negli studi in cui chi somministra il trattamento non è in cieco si pone comunque il nascondere la sequenza di allocazione fino a che non è effettivamente utilizzata

**Definition 13.5.2** (Allocation concealment). Mantenimento della segretezza della lista di randomizzazione sino alla richiesta di randomizzazione, per evitare allocazioni non casuali (es pazienti più gravi o amici/parenti)

*Remark 31.* Per garantirla la lista non deve essere in possesso degli sperimentatori, ma di soggetti terzi alla sperimentazione (ma di supporto alla stessa) oppure in sistemi automatici.

### 13.5.2 Tipologie

**Definition 13.5.3** (Semplice). Quella che si ottiene lanciando una moneta

*Remark 32.* Semplice da effettuare ma non vi è garanzia su

- bilanciamento complessivo dei bracci a fine studio: sbilanciamenti notevoli su
- bilanciamento in itinere durante lo studio: ad esempio si evita che grossa parte dei trattati (o dei controlli) vengano trattati all'inizio dello studio e i controlli alla fine (o viceversa), esponendo magari a differenti condizioni ambientali le due popolazioni
- bilanciamento tra trattati e controlli su fattori prognostici correlati dell'outcome

**Definition 13.5.4** (A blocchi bilanciati). Non si randomizza un singolo paziente ma un blocco bilanciato di pazienti (es in ratio 1:1, due o 4 o 6 pazienti equamente suddivisi tra i due bracci di trattamento)

*Remark 33.* È good practice effettuare blocchi di dimensione variabile, per aumentare l'allocation concealment.

*Remark 34.* Rispetto alla randomizzazione semplice il blocking risolve lo sbilanciamento tra bracci complessivo (o quello in itinere durante lo studio)

**Definition 13.5.5** (Stratificata). Realizzazione di una lista di randomizzazione (semplice o a blocchi) per ogni strato definito da un fattore prognostico/predittivo

*Remark 35.* Quest'ultima garantisce il fatto che non vi siano sbilanciamenti grossi sui fattori prognostici/predittivi tra casi e controlli.

Nei multicentrici d'obbligo è farla utilizzando il centro come fattore di stratificazione: sia per evitare che tutti i trattamenti siano fatti in un centro e tutti i controlli in un altro, sia perché eventuali problemi con un centro che portino all'esclusione dello stesso non hanno ripercussioni sul bilanciamento complessivo

*Remark 36.* La considerazione di molteplici fattori di stratificazione contemporaneamente pone il rischio di sbilanciamenti numerici tra trattati e controlli, a maggior ragione se le liste non sono costruite mediante blocchi bilanciati.

Soprattutto su trial di piccole dimensioni questo conduce alla minimizzazione (o randomizzazione adattiva)

**Definition 13.5.6** (Adattata (minimization)). L'idea è, ad ogni nuovo paziente da reclutare, tenere conto delle sue caratteristiche per valutare quale assegnazione introduce minor sbilanciamento; l'assegnazione poi può avvenire in maniera deterministica o probabilistica (es assegnando

### 13.5.3 Altre questioni

**Comparabilità dei gruppi randomizzati** [1] sostiene come:

- a causa della randomizzazione non sia necessario valutare la comparabilità dei gruppi randomizzati (la  $p$  sarebbe la probabilità della differenza osservata tra gruppi nell'ipotesi non vi fosse differenza tra gruppi; ma effettivamente non vi è differenza tra gruppi nella popolazione, essendo avvenuta la randomizzazione)
- può essere necessario aggiustare per sbilanciamenti, soprattutto se le variabili che in seguito alla randomizzazione risultino sbilanciate abbiano una associazione con l'outcome (es età). Il modo per procedere è aggiornare le stime di efficacia per fattori che per sensibilità clinica e buon senso appaiono sbilanciati.

## 13.6 Definizione dell'effetto del trattamento

Per definire univocamente il *segnale*, ovvero la grandezza attraverso la quale viene valutato, a livello di gruppo e in termini comparativi l'effetto del trattamento sperimentale che si intende studiare occorre effettuare una serie di passaggi [4].

1. Definizione degli aspetti della malattia in studio su cui il trattamento vuole incidere; questi sono definiti come *livelli terapeutici*
2. Identificazione di una sola variabile (quindi di un dato processo di misurazione) utile ad identificare l'effetto del trattamento su ogni soggetto,

per ogni livello terapeutico considerato. Questa variabile è chiamata **end-point**, o variabile risposta. La sua scelta dovrebbe esser determinata primariamente da motivi clinici.

3. Definizione degli **indicatori di gruppo**: decidere come sintetizzare ciascun end point a livello di gruppo; la scelta di questo dipende dal tipo di end-point e dalla sua distribuzione (la scelta qui è primariamente statistica)
4. Per ciascun indicatore di gruppo bisogna definire il modo in cui vengono comparati matematicamente i gruppi a confronto (differenza o rapporto tra indicatori di gruppo). Questo è il **segnale** ovvero l'effetto complessivo del trattamento
5. Occorre **gerarchizzare i livelli terapeutici**, identificando di conseguenza il segnale primario, sulla quale si giudicherà il confronto dei trattamenti in studio
6. Definizione delle **soglie di rilevanza/non rilevanza** clinica. Nel caso di *studi di superiorità* (volti a dimostrare la superiorità del trattamento sperimentale) bisogna formulare una previsione dell'entità dell'effetto del gruppo sperimentale considerata sufficiente per dichiarare l'effetto biologicamente/clinicamente rilevante. Negli studi di equivalenza e non inferiorità bisogna definire di non rilevanza clinica o margine di equivalenza, cioè la massima differenza che si può tollerare tra gruppi per poter affermare che sono simili (studi di equivalenza) o che quello sperimentale sia non inferiore a quello di controllo (non inferiorità).

## 13.7 Disegni meno frequenti

*Remark 37.* Tipicamente i trial prevedono due bracci paralleli in allocazione 1:1; questo garantisce potenza maggiore per un confronto su una singola ipotesi. Ciò non toglie che in determinate circostanze altri disegni possano tornare comodi o esser necessari.

### 13.7.1 Disegno a più bracci paralleli

*Remark 38.* Essendo molteplici le ipotesi da saggiare sono necessari più test statistici, tra loro non indipendenti; occorre pertanto gestire problemi legati alla numerosità campionaria e molteplicità.

### 13.7.2 Disegno fattoriale

*Remark 39.* Mediante questo si effettuano 2 o più confronti terapeutici diversi nella stessa ricerca, senza aumentare il numero totale di pazienti

*Remark 40.* L'assunto che sta alla base di questo disegno è che le terapie abbiano meccanismi d'azione diversi e che non interferiscano tra loro (ossia effetti additivi, senza interazione); se viceversa si prevede di saggiare l'interazione è necessario dimensionare adeguatamente (cosa che solitamente porta a numerosità elevate)



*Remark 41.* Ad esempio in un fattoriale  $2 \times 2$  (due trattamenti per entrambi di due livelli, presente o assente) la lista di randomizzazione contempererà quattro bracci.

### 13.7.3 Disegno cross-over

*Remark 42.* In questo disegno ogni partecipante è controllo di se stesso; nella forma più semplice il disegno prevede che ogni paziente riceva entrambi i trattamenti a confronto in successione, secondo un ordine casuale

*Remark 43* (Condizioni di applicabilità). Quando:

- le terapie sono di durata breve
- gli effetti si manifestano in poco
- in assenza di trattamento la malattia rimane stabile nel periodo necessario per somministrare i due trattamenti
- gli effetti a lungo termine del primo trattamento sono nulli o comunque scompaiono dopo un periodo di *wash-out* adeguato

*Remark 44.* Trova scarsa applicazione in oncologia: può esser applicato nel campo delle terapie palliative (es terapia del dolore per confrontare due diversi analgesici)

## 13.8 Altri studi comparativi (metodologicamente inferiori)

### 13.8.1 Prima-dopo

**Definition 13.8.1** (Disegno). Ad un gruppo di pazienti somministriamo il trattamento e confrontiamo il risultato dopo con quello prima

*Remark 45* (Pro). Sono:

- Disegno semplice (conduzione/comunicazione al pz.)
- Minori problemi etici se vi è presunzione di efficacia (tutti trattati, nessun escluso)
- Disegno parsimonioso (minor variabilità, test paired più potenti, sufficienza di un campione più contenuto rispetto ai disegni controllati, minor costo dello studio)

*Remark 46* (Contro). Rischio di introduzione di **bias** (la stima di effetto del trattamento è sbagliata) in presenza di:

1. **Dinamica spontanea** della malattia: negli studi prima dopo si assume che in assenza dell'intervento non ci sarebbe stata una variazione della malattia. È un assunto non verificabile; si possono fare congetture (che dipendono dalla patologia in questione) ma rimangono tali.  
La direzione del Bias dipende dalla direzione della dinamica: es se la malattia “migliora” nel tempo e il trattamento mostra un effetto positivo,

l'efficacia del trattamento in una stima prima- dopo sarà “gonfiato”)  
 Il problema è tanto maggiore (e il disegno prima-dopo è meno applicabile) tanto quanto: La malattia presenta una dinamica spontanea veloce e il trattamento è lento a produrre effetti

2. **Variazioni del contesto** dello studio: le variazioni tra il prima e il dopo che incidono sull'outcome (e/o sulla sua valutazione) oltre al trattamento. In primis: Cambio Outcome Assessor, Macchinari, Software di valutazione.  
 Rischio più alto tanto è più lungo lo studio; difficile prevedere la direzione del bias.
3. **Regressione verso la media:**
4. **Effetto apprendimento:** se la valutazione dell'effetto del trattamento avviene mediante test di “abilità” del paziente, questi può imparare e performare meglio al test anche in assenza di efficacia del trattamento. (Una soluzione parziale potrebbe essere una fase di pratica pre valutazione basale)
5. **Effetto psicologico** (placebo): solo il fatto di saper di esser curati, ha effetti benefici sulla prognosi della malattia, indipendentemente dall'efficacia clinica del trattamento.  
 L'effetto placebo tende ad incrementare la stima di efficacia di un trattamento in un confronto pre-post; a contrario del disegno con gruppo di controllo trattato a placebo, in un pre-post questo bias non riesce ad esser scorporato.

### 13.8.2 Trial controllati non randomizzati

*Remark 47.* Rientrano in questa categoria studi eterogenei con

- Controlli paralleli, non selezionati mediante il caso (es iniziale del cognome, giorni sett.)
- Controlli storici o da banche dati (pz trattati in passato)

*Remark 48.* Rimane dubbia la confrontabilità tra i due gruppi: pertanto considerati da diversi autori/istituzioni metodologicamente di qualità inferiore!

*Remark 49.* Ergo: se possibile non progettarli. Nel caso, l'analisi non può prescindere da modelli multivariati.

# Capitolo 14

## Fase 4

Sono gli studi sul farmaco, spesso osservazionali, svolti nell'ambito della pratica clinica dopo che questo è stato messo in commercio.

**Razionale** Negli studi di fase precedente la numerosità di pazienti nonché il loro follow up è relativamente limitato, pertanto di fatto si riescono a rilevare solamente gli eventi avversi più comuni e che si manifestano entro relativamente poco tempo.

**Obiettivi** Si pongono la valutazione di [amadori2004sperimentazioneclinicaoncologia ]:

- efficacia e tollerabilità del trattamento in popolazioni non selezionate
- interazioni con altri farmaci in commercio
- impiego a lungo termine del farmaco (es terapie croniche)
- implicazioni farmacoeconomiche
- impatto sulla qualità della vita

Alla base di diversi obiettivi vi è l'attività di *farmacosorveglianza/farmacovigilanza*, che si basa sulle segnalazioni spontanee delle reazioni avverse, note o impreviste

**Disegno** Tipicamente si tratta di disegni osservazionali, quindi coorte, caso-controllo o studi cross section. Non servono a dimostrare l'efficacia di un trattamento (nota al termine della fase III) ne hanno il rigore metodologico degli studi sperimentali nel dimostrare associazioni causali.



Parte V

Studi osservazionali



# Capitolo 15

## Introduzione osservazionali

### 15.1 Tipi di studi

**Definition 15.1.1** (Coorte). Il ricercatore seleziona due gruppi di soggetti, con e senza la caratteristica in studio (es genere: maschi e femmine) ma senza l'evento di interesse (es tumore). I due gruppi vengono osservati per un dato periodo di tempo e viene confrontata l'incidenza dell'evento di interesse.

*Remark 50.* A parte la caratteristica i due gruppi dovrebbero essere il più possibile omogenei: quindi ben venga prelevare dallo stesso bacino e laddove fattibile formare il campione mediante estrazione casuale dalle due popolazioni (con/senza fattore).

**Definition 15.1.2** (Caso controllo). Il ricercatore seleziona due gruppi di soggetti, rispettivamente con (casi) e senza (controlli) l'evento di interesse oggi e per ogni soggetto ricerca nel passato, per un dato periodo di osservazione, informazioni sull'esposizione alla caratteristica studiata. Si stima/confronta l'associazione tra esposizione ed esito.

*Remark 51.* Alcune considerazioni:

- anche in questo caso si può estrarre casualmente tra casi e controlli per avere rappresentatività delle popolazioni
- la scelta dei controlli è particolarmente critica

**Definition 15.1.3** (Cross sectional). Si rileva in un dato momento la presenza di un determinata condizione

*Remark 52.* Questi studi non sono adatti per inferire causalità in quanto nel caso di esposizione ed esito con una fotografia in un dato momento non si può stabilire cosa è venuto temporalmente/clinicamente prima.

### 15.2 Bias e confondimento

Di fronte ad un risultato che mostri un'associazione tra esposizione e malattia, il ricercatore si deve chiedere:

- può esser dovuta al caso? Nell'analisi statistica e nello specifico nella determinazione dell'ampiezza campionaria si stabilisce come gestire questo tipo di errore (di prima specie)
- può esser dovuta a **bias**/distorsione?
- può esser dovuta a **confondimento**?

**Definition 15.2.1** (Bias). Qualunque errore sistematico che porti ad una stima non corretta dell'associazione tra esposizione ed evento

*Remark 53.* Può verificarsi in qualunque fase dello studio e non riesce ad esser eliminato in sede di analisi statistica

*Remark 54.* Alcuni esempi sono

- distorsione da *selezione*: l'errore sistematico riguarda la selezione dei soggetti da includere nello studio
- distorsione di *osservazione*: difficoltà a ricordare, influsso dell'intervistatore, perdita di soggetti al follow up sistematicamente diversi da coloro che rimangono,

**Definition 15.2.2** (Confondimento). Si ha se l'entità e a volte la direzione dell'associazione tra caratteristica ed evento è modificata dalla presenza di un terzo fattore che è contemporaneamente:

- associato con la caratteristica
- distribuito in modo sbilanciato tra i gruppi a confronto

*Remark 55.* Nel bias lo sbilanciamento sistematico tra i gruppi è indotto dal disegno, mentre nel confondimento è insito nel fenomeno.

*Remark 56.* A volte il termine confondimento è usato impropriamente come sinonimo di associazione spuria, che dai medici è intesa come associazione “non vera” mentre dagli statistici per indicare una associazione tra due fattori in realtà causata da un terzo (es correlazione tra consumo di cioccolata a livello nazionale e premi nobel è una associazione spuria per la presenza del PIL procapite)



# Capitolo 16

## Coorte

In uno studio di coorte gli individui sono seguiti per un dato periodo di tempo per monitorare il loro stato di salute e nello specifico, uno specifico evento (es morte, ammalarsi ecc).

### 16.1 Disegni

**Disegno base** L'approccio più semplice consiste nel selezionare due gruppi di persone al baseline; il primo consiste in persone che possiedono un qualche attributo di interesse (es esposizione ad un fattore di rischio) mentre l'altro no, al fine di studiare l'eccesso di mortalità/morbidità associato al fattore di rischio. A volte si riesce a studiare la popolazione completa esposta al fattore di rischio; quando invece se ne può studiare solo un campione è meglio sceglierlo in maniera casuale.

Ogni soggetto che entra nello studio deve essere libero da evento analizzato (morte/malattia) al baseline proprio perché si sta cercando di capire la causalità associata al fattore di rischio.

**Disegno con campione unico** Prendiamo tot pazienti in blocco unico ed analizziamo l'associazione in questi:

- pro l'informazione sulla stratificazione in pazienti con e senza non è necessaria a priori e in generale è logisticamente più facile da fare;
- pro: è un disegno comodo soprattutto se con un unico studio si stanno analizzando il contributo di più fattori di rischio contemporaneamente;
- pro: incidentalmente si possono fornire stime della prevalenza del fattore di rischio (se si prende un campione casuale della popolazione);
- contro: la distribuzione del fattore di rischio potrebbe essere notevolmente e originare confronti statistici poco potenti.

**Disegno con confronto esterno** Arruolare non due bracci ma solo il braccio di coloro che hanno l'esposizione di interesse e confrontare l'esito con una popolazione esterna, spesso e volentieri la popolazione generale. Costa meno

ma la possibilità di bias è indiscussa; la popolazione generale può differire dal campione non solo per il fattore considerato ma per molto altro che si fatica a controllare.

**Coorte retrospettiva** Magari svolto su dati già disponibili, è comunque necessario:

- garantire che all’inizio del follow up tutti i pazienti fossero disease free
- essere confidenti sul fatto che il dataset collezionato per altri fini possa rispondere anche a quello dell’indagine (es valutare definizioni adottate ed eventuali cambiamenti nel tempo, completezza dei dati eccetera)

## 16.2 Pro/contro

Vantaggi dello studio di coorte:

- la prospettività, unita al fatto che i pazienti siano liberi da evento all’inizio del follow up, è ideale per suggerire la causalità delle associazioni
- l’effetto di un singolo fattore di rischio può essere valutato su più malattie (outcome)

Svantaggi:

- costosi (es se l’outcome necessita di esami) e lunghi (in particolar modo per malattie, outcome, lenti a svilupparsi) se condotti in maniera prospettica
- non adatti per outcome rari: in tal caso potremmo dover arruolare un campione enorme (al fine che si verifichino eventi) o monitorare per tempi molto lunghi. Per questo si possono scegliere outcome intermedi (surrogati) più frequenti accettandone tutte le critiche associate (validazione, non è un outcome hard eccetera);
- study effect (assenza di blinding): se il paziente sa di essere in uno studio potrebbe comportarsi diversamente dal caso in cui non lo sapesse. E le associazioni derivanti potrebbero risultare falsate o non rappresentative della realtà
- l’esposizione al fattore di interesse potrebbe variare nel corso dello studio (può essere considerato in analisi)
- pazienti che si ritirano o persi al fup: gestibili a patto che il censoring non sia informativo

## 16.3 Strategie d’analisi

Se ciascuno nella coorte viene arruolato nello stesso momento e viene seguito per lo stesso ammontare di tempo si possono analizzare i dati in maniera binaria classica (tipicamente). Questa è una **fixed cohort**.

Nel caso il follow up non sia uguale per tutti (**variable cohort**) meglio approcciarsi con survival analysis oppure analisi person-years (poisson)

## Capitolo 17

### Caso controllo



Parte VI

Studi di diagnostica



## Capitolo 18

# Introduzione all diagnostica

### 18.1 Introduzione

La diagnostica è volta a raccogliere ed analizzare informazioni per determinare la *condizione del paziente*. La ricerca diagnostica è motivata dal fatto che non sempre un metodo di diagnosi efficace esista; altresì si può esser interessati a sviluppare nuove metodologie diagnostiche se quelle esistenti sono particolarmente costose, invasive, tossiche ecc.

Dal punto di vista statistico, la metodologia della diagnostica, volta prevalentemente alla classificazione, si applica anche ad altri problemi in medicina; ad esempio anche la prognosi può esser considerata una diagnosi (si cerca di indovinare l'esito, non la diagnosi). Allo stesso modo i test di screening sono test diagnostici (in popolazioni con prevalenze tipicamente rare).<sup>1</sup>

La performance di un test diagnostico può esser valutata a livelli progressivi; qualità del dato/immagine (aspetto tecnico), accuratezza diagnostica, effetto sulle decisioni di trattamento, impatto sull'outcome del paziente, costo/efficacia per la società (aspetti scientifici). In ogni modo per esser efficace a livelli superiori, deve esser efficace ad un livello inferiore. Qui studiamo il primo degli aspetti scientifici che si pongono ovvero lo studio dell'accuratezza diagnostica del test.

Lo sviluppo di un nuovo test, seguendo [21], è utile in presenza delle seguenti condizioni (riguardanti la malattia, il suo trattamento, e il test stesso):

1. la malattia dovrebbe esser grave o potenzialmente grave (se il diagnosticarla non provoca un risparmio di quantità/qualità di vita, non è costo efficace)
2. la malattia dovrebbe esser abbastanza prevalente nella popolazione target (non rarissima)
3. la malattia dovrebbe esser trattabile (se no è inutile testare)
4. il trattamento dovrebbe esser disponibile per coloro che hanno un test positivo (es non troppo costoso, se no è comunque inutile)

---

<sup>1</sup>Si tratta di studi in cui si desidera, per diversi fini, classificare i pazienti; che la classificazione riguardi uno stato di malattia o un esito, dal punto di vista strettamente di calcolo è lo stesso.

5. il test non dovrebbe causare (troppi) effetti avversi al paziente
6. il test dovrebbe classificare accuratamente gli individui malati da quelli sani: questa è l'*accuratezza diagnostica*

### 18.1.1 Disegni di ricerca principali

Sebbene nel seguito ci si concentri sullo studio di accuratezza, esistono differenti tipologie di studi nell'ambito della diagnostica ed essi dipendono sostanzialmente dall'obiettivo/domanda della ricerca.

Se l'obiettivo è valutare la *capacità discriminativa di un test nell'individuare sani e malati*, lo studio principale è quello di **accuratezza diagnostica** in cui viene indagata la relazione tra risultato del test e presenza della condizione. Hanno tipicamente un disegno cross-sectional.

Se si vuole valutare l'*impatto dell'impiego del metodo diagnostico nella pratica clinica e sulla prognosi* dei pazienti, esistono diversi disegni, tra i quali l'**RTC diagnostico** è metodologicamente quello più forte. Anche disegni di coorte (guardo chi è stato diagnosticato con che cosa, e ne seguo poi il fup) e caso controllo (guardo malati e non malati oggi e torno indietro nel ricercare procedure diagnostiche di mio interesse) sono impiegati

La *sintesi delle evidenze disponibili* può esser affrontata mediante **revisioni sistematiche** (che forniscono una valutazione globale della procedura diagnostica).

Sempre in questa categoria, esistono altresì **studi di economia sanitaria** volti alla costo efficacia delle procedure diagnostiche. Vengono infine svolti studi sullo sviluppo di **equazioni di predizione** (o *clinical prediction rules, CPR*) diagnostica/prognostica, che possano servire al clinico per effettuare delle decisioni nella propria pratica.

### 18.1.2 Architettura della ricerca diagnostica

Qui mettere buntinx pag 20

## 18.2 RCT Diagnostici

Knotterus capitolo 4, TODO



## Capitolo 19

# Studi di accuratezza diagnostica

### 19.1 Introduzione

Gli **studi di accuratezza diagnostica**, che costituiscono il nostro focus, sono ricerche volte a determinare la *capacità di un test di discriminare* tra pazienti con e senza una determinata condizione clinica.

Qualunque dato (caratteristiche del paziente, segni e sintomi, esami fisici, storia del pz oppure test di laboratorio) può esser in linea teorica considerato come “test”, e verificarne la capacità discriminatoria.

Sostanzialmente un test accurato classifica i soggetti correttamente in relazione alla loro condizione. Test inaccurati fanno sì che (troppi) individui malati vengano classificati come sani e viceversa. I primi non vengono trattati come dovrebbero, i secondi possono ricevere procedure non necessario, spesso invasive e costose.

Pertanto prima che un test possa esser utilizzato in pratica clinica, la sua accuratezza diagnostica deve esser valutata.

Per poter valutare la capacità classificatoria di un test, i suoi risultati debbono esser confrontati alternativamente con una diagnosi standard di riferimento; un **gold standard**, ovvero un altro test che classifichi in assoluta correttezza la condizione del pz non necessariamente esiste<sup>1</sup>. Nel caso, occorre impiegare un *reference standard* che approssimi il gold standard al meglio.

In merito alla presenza di un gold standard, va comunque precisato che gli studi di accuratezza diagnostica non sono meramente studi di “agreement” tra due misure. Interpretare ogni differenza tra lo standard e il test sperimentale come un fallimento di quest’ultimo non è necessariamente corretto (soprattutto se lo standard non è gold) . Può esser peraltro che due metodi misurino concetti lievemente differenti.

---

<sup>1</sup>Esempi di gold standard sono i report dall’autopsia, i rilievi chirurgici, risultati dell’analisi patologica su campioni

Test (Y)	Reference standard (D)	
	Presente (1)	Assente (0)
Presente (1)	$s_1 = Tp$	$r_1 = Fp$
Assente (0)	$s_0 = Fn$	$r_0 = Tn$

Tabella 19.1: Matrice 2x2

## 19.2 Misure di accuratezza diagnostica

### 19.2.1 Dati dicotomici

Se ci riferiamo ad un test che fornisce nativamente risultati in forma dicotomica (del tipo malattia assente o presente) e confrontando i suoi risultati con un gold o reference standard. Arriviamo a definire la tabella 19.1, e i valori (conteggi) di cui è composta:

- **Tp**: true positive. Casi in cui il test *individua* ( $Y=1$ ) *correttamente* ( $D=1$ ) la presenza di malattia (effettivamente *presente*).
- **Fp**: false positive. Casi in cui il test *individua* ( $Y=1$ ) *fallacemente* ( $D=0$ ) la presenza di malattia (effettivamente *assente*).
- **Fn**: false negative. Casi in cui il test *esclude* ( $Y=0$ ) *erroneamente* ( $D=1$ ) la malattia (effettivamente *presente*).
- **Tn**: true negative. Casi in cui il test *esclude* ( $Y=0$ ) *correttamente* ( $D=0$ ) la malattia (effettivamente *assente*).

#### 19.2.1.1 Sensibilità, specificità

La **sensibilità** (*sensitivity* o *true positive fraction*, TPF) è l'abilità del test di individuare la malattia quando è presente:

$$Sens = TPF = P(Y = 1|D = 1) = \frac{Tp}{Tp + Fn} \quad (19.1)$$

La **specificità** (*specificity* o *true negative fraction*, TNF) è l'abilità del test di escludere la condizione in pazienti che non ne siano affetti.

$$Spec = TNF = P(Y = 0|D = 0) = \frac{Tn}{Tn + Fp} \quad (19.2)$$

Un test perfetto (che non commette falsi positivi o negativi) avrebbe:

$$Sens = Spec = 1 \quad (19.3)$$

Un test non in grado di discriminare (equivalente a tirare una moneta sia nel caso il paziente abbia o non abbia la malattia):

$$Sens = Spec = 0.5 \quad (19.4)$$

Allo stesso modo possono esser definite le seguenti misure:

- FNF, o *false negative fraction*, come

$$FNF = P(Y = 0|D = 1) = \frac{Fn}{Tp + Fn} = 1 - TPF \quad (19.5)$$

- FPF, o *false positive fraction*

$$FPF = P(Y = 1|D = 0) = \frac{Fp}{Tn + Fp} = 1 - TNF \quad (19.6)$$

La **prevalenza** della malattia nella popolazione è definita come:

$$Prev = \rho = P(D = 1) = \frac{Fn + Tp}{n} \quad (19.7)$$

dove  $n = Fn + Tp + Tn + Fp$ .

L'**accuratezza** è definita come la probabilità che il test azzechi la diagnosi;

$$Acc = P(Y = D) = \frac{Tp + Tn}{n} \quad (19.8)$$

La **probabilità di missclassification** è il complemento a 1 dell'accuratezza e può esser scritta come funzione della prevalenza della malattia, di FNF e FPF

$$Miss = P(Y \neq D) = \rho(FNF) + (1 - \rho)FPF \quad (19.9)$$

Le misure complessive di accuratezza/missclassification non sono generalmente considerate una sintesi adeguata dell'accuratezza diagnostica di un test medico. Piuttosto, bisognerebbe riportare sia sensibilità e specificità (ovvero i complementi a 1 di FNF e FPF) separatamente poiché:

- costi e conseguenze dei due tipi di errori possono esser molto differenti (falsi negativi perdono le cure necessarie, falsi positivi si sottopongono a terapie non necessarie). Solitamente i falsi positivi sono giudicati meno gravi dei falsi negativi;
- inoltre mentre sensibilità e specificità non dipendono dalla prevalenza della malattia nella popolazione, misure di sintesi come accuratezza e missclassification vi dipendono largamente.

### 19.2.1.2 Valori predittivi

Guardando alla tabella per riga, come alternativa alle misure di accuratezza che fondano il proprio denominatore sullo stato della malattia vi sono quelle che lo fondano sul risultato del test, che pongono enfasi su quanto bene i risultati del test prevedano lo stato effettivo di malattia.

Il valore predittivo positivo (**positive predictive value** o PPV) esprime la probabilità che un test positivo riesca ad individuare correttamente un soggetto avente la condizione:

$$PPV = P(D = 1|Y = 1) = \frac{Tp}{Tp + Fp} \quad (19.10)$$

Il valore predittivo negativo (**negative predictive value** o NPV) esprime la probabilità che un test negativo riesca ad escludere correttamente che un soggetto sia senza la condizione:

$$NPV = P(D = 0|Y = 0) = \frac{Tn}{Tn + Fn} \quad (19.11)$$

Un test perfetto prevederà la presenza di condizione in maniera perfetta, avendo

$$PPV = NPV = 1 \quad (19.12)$$

Invece un test inutile che non porta informazione sulla presenza effettiva di malattia sarà tale che

$$PPV = P(D = 1|Y = 1) = P(D = 1) = \rho \quad (19.13)$$

$$NPV = P(D = 0|Y = 0) = P(D = 0) = 1 - \rho \quad (19.14)$$

I valori predittivi non sono usati per quantificare la performance intrinseca del test, perchè non hanno come denominatore la condizione effettiva e dipendono tra l'altro dalla prevalenza della malattia. Piuttosto sono impiegati per quantificare il valore clinico del test (a questi sono maggiormente interessati paziente e caregiver).

In generale i valori predittivi dipendono sia dalla performance intrinseca del test (sensibilità, specificità) che dalla prevalenza della malattia. Si può scrivere

$$PPV = \frac{\rho \cdot Sens}{\rho \cdot Sens + (1 - \rho) \cdot (1 - Spec)} \quad (19.15)$$

$$NPV = \frac{(1 - \rho) \cdot Spec}{(1 - \rho) \cdot Spec + \rho \cdot (1 - Sens)} \quad (19.16)$$

La dimostrazione è un'applicazione del teorema di Bayes.

### 19.2.1.3 Uso di R - Stima accuratezza

Esempio 2.1 pag 17 della pepe

```
cass <- read.csv("~/dataset/pepe/est1.csv")

## Warning in file(file, "rt"): non è possibile aprire il file '/home/l/dataset/pepe/
File o directory non esistente
## Error in file(file, "rt"): non è possibile aprire la connessione

names(cass) <- c("cad", "est", "cph")

## Error: oggetto 'cass' non trovato

## Cad: Coronary artery disease
## EST: Exercise Stress test
## CPH: Chest pain history
dim(cass)

## Error: oggetto 'cass' non trovato
```

```
head(cass)

## Error:  oggetto 'cass' non trovato

library(lbdiag)

## Error in library(lbdiag):  non c'è alcun pacchetto chiamato 'lbdiag'

da(test=cass$est, refstd=cass$cad)

## Error in da(test = cass$est, refstd = cass$cad):  non trovo la funzione
"da"
```

Il commento si può fare per colonne o per righe della tabella:

- per “colonne”, la prevalenza della malattia (prev) è molto alta, quasi il 70%. La sensibilità è circa dell’80 mentre la specificità è del 74%. Il test manca il 20% dei malati e erroneamente identifica come malati il 26% dei sani. La decisione di utilizzare il test necessiterebbe di prendere in considerazione il rischio e il beneficio di procedure diagnostiche aggiuntive e/o trattamenti associati ad una diagnosi positiva, cosiccome le conseguenze delle mancate diagnosi.
- per “righe”, circa il 63% della popolazione viene diagnosticata positiva al test. Tra i soggetti che hanno test positivo la stragrande maggioranza (88%) ha effettivamente la malattia. I positivi quindi sono verosimilmente all’ultima fase diagnostica dato che hanno una elevata probabilità di esser malati  
Al contrario il 39% dei soggetti che vengono diagnosticati negativi hanno in realtà la malattia. Pertanto potrebbe esser di interesse la ricerca di ulteriori tecniche diagnostiche da applicare a coloro che hanno un test negativo, al fine di identificare coloro che in realtà hanno la malattia.

#### 19.2.1.4 Uso di R - Inferenza accuratezza

Esempio 2.3 pag 22 della pepe: intervalli di confidenza delle stime

```
## Sens 95CI
binom.test(815,1023)$conf.int

## [1] 0.7706868 0.8209461
## attr(,"conf.level")
## [1] 0.95

## Spec 95CI
binom.test(327,442)$conf.int

## [1] 0.6962660 0.7801285
## attr(,"conf.level")
## [1] 0.95

## PPV 95CI
binom.test(815,930)$conf.int
```

```
## [1] 0.8534513 0.8968204
## attr(,"conf.level")
## [1] 0.95

## NPV 95CI
binom.test(327,535)$conf.int

## [1] 0.5684514 0.6527429
## attr(,"conf.level")
## [1] 0.95

## Joint 95 Confidence interval Sens spec pag 23
binom.test(815,1023, conf.level = 0.975)$conf.int

## [1] 0.7669239 0.8242471
## attr(,"conf.level")
## [1] 0.975

binom.test(327,442, conf.level = 0.975)$conf.int

## [1] 0.6900029 0.7855267
## attr(,"conf.level")
## [1] 0.975
```

#### 19.2.1.5 Molteplicità di focus diagnostici entro paziente

Nel caso in cui la **presenza di malattia possa esser diagnosticata più volte per ogni paziente** (es se si vuole diagnosticare se una lesione è maligna, ma le lesioni in un pz possano esser molteplici), allora si può costruire la tabella a livello di lesione/polipo, anzichè al livello di paziente. Per farlo i conti vanno fatti a livello di singola lesione, per cui un Tp sarebbe una lesione che il test ha correttamente individuato come positiva.

Nel caso invece si desidera di *passare da livello di lesione a quello di paziente*:

- per il pz i-esimo il reference standard è impostato a malato se il pz stesso ha *almeno una lesione maligna*
- per il pz i-esimo il reference standard è impostato a sano se il pz stesso non ha *neanche una lesione maligna*

#### 19.2.2 Dati quantitativi

Molti test forniscono una misura numerica.

In questo setting, il valore restituito dal test varia da 0.03 a 0.58 nei pazienti con frattura, da 0 a 0.13 nei pazienti senza frattura.

Siamo interessati a diagnosticare la frattura sulla base del test numerico, meno invasivo

```
gap <- c(0.58, .41, .18, .15, .15, .10, .07, .07, .05, .03,
        .13, .13, .07, .05, .03, .03, .03, 0, 0, 0)
```

```
fracture <- c(rep(1,10),rep(0,10))
hv <- data.frame(fracture,gap)
hv[hv$fracture==1 , "gap"]

## [1] 0.58 0.41 0.18 0.15 0.15 0.10 0.07 0.07 0.05 0.03

hv[hv$fracture==0 , "gap"]

## [1] 0.13 0.13 0.07 0.05 0.03 0.03 0.03 0.00 0.00 0.00
```

Per poter determinare sensibilità e specificità dobbiamo scegliere una **soglia** e *specificare se sopra di essa il test sia considerato positivo o negativo*. Nel nostro caso, data le conoscenze cliniche, specifichiamo che un valore di test superiore a 0.05 (decision threshold) sia considerabile come un test positivo (quindi predittivo della presenza di frattura).

Avremmo potuto scegliere qualsiasi valore come soglia per determinare l'esito del test. Questo fa sì che a differenti soglie corrispondano differenti performance diagnostiche

### 19.2.3 Dati ordinali





## Parte VII

# Revisioni sistematiche



# Capitolo 20

## Introduzione

### 20.1 Definizioni e risorse utili

*Remark 57.* Qui ci si basa per lo più su [18]

**Definition 20.1.1** (Revisione sistematica (RS)). Studio che sintetizza l'evidenza empirica rispettante pre-determinati criteri di eleggibilità, al fine di rispondere ad una specifica domanda di ricerca

*Remark 58.* Le RS si contrappongono alle Revisioni Narrative

**Definition 20.1.2** (Revisione narrativa). Revisione dove l'autore non esplicita come ha scelto gli studi da includere o meno (scegliendo mediante esperienza personale o contatti/conoscenze)

*Remark 59* (Caratteristiche di una revisione sistematica). Rispettivamente:

- obiettivo/set di obiettivi chiari, con criteri di eleggibilità degli studi pre-determinati
- metodologia riproducibile
- ricerca sistematica che identifichi tutti gli studi che rispettino i criteri di eleggibilità
- valutazione validità dei risultati degli studi selezionati (es valutazione rischio di bias)
- sintesi sistematica di caratteristiche e risultati degli studi inclusi

*Remark 60* (Fasi di una revisione sistematica). Rispettivamente:

- Dichiarare gli obiettivi e i criteri di eleggibilità
- Ricerca di studi che sembrano rispettare i criteri di eleggibilità
- Porre in dataset le caratteristiche degli studi identificati e valutarne la qualità metodologica
- applicare i criteri di eleggibilità e giustificare ogni esclusione

- raccogliere i dati da analizzare
- analisi dei risultati degli studi eleggibili, usando meta-analisi se appropriato e possibile
- analisi di sensibilità e analisi sottogruppi, se appropriato e possibile
- preparazione del report

**Definition 20.1.3** (Metanalisi (MA)). Metodo statistico di sintesi dei risultati di differenti studi.

*Remark 61.* Non necessariamente tutte le revisioni sistematiche presentano anche una metanalisi: può essere talvolta non utile o appropriato combinare quantitativamente le informazioni derivanti da studi fra loro troppo diversi o eterogenei.

*Remark 62.* L'analisi della coerenza e qualità di un insieme di studi è una delle caratteristiche più importanti di una RS

*Remark 63.* Si può ritenere appropriata la combinazione quantitativa dei dati di studi diversi (MA) quando:

- più di uno studio ha stimato l'effetto del trattamento/terapia
- le differenze fra gli studi in termini di pazienti, interventi e caratteristiche del setting sono minimo o comunque non permettono a priori, di ipotizzare un impatto sull'outcome
- l'outcome nei diversi studi è stato misurato in maniera simile
- gli autori degli studi primari riportano i dati numerici necessari per effettuare la combinazione

*Remark 64* (Risorse per il lettore interessato). In alcun ordine in particolare:

- Cochrane Handbook for Systematic Reviews of Interventions: manuale per revisioni su studi di efficacia
- Cochrane Handbook for Systematic Reviews of Diagnostic Test Accuracy: manuale per studi di accuratezza diagnostica
- [liberati2009prisma]: guideline scrittura revisione sistematica e metanalisi
- RevMan e RevMan Web: software per la conduzione di revisioni della Cochrane collaboration
- Altri tool cochrane

## 20.2 Preparazione e mantenimento di una review (Cochrane)

### 20.2.1 Protocollo

Dato che lo svolgimento di una revisione necessita di diverse decisioni, e volendo evitare che queste siano indirizzate dai risultati degli studi inclusi, analogamente

a quanto avviene per i singoli studi (es rct) la pubblicazione del protocollo (anch'esso fattibile mediante revman e pubblicabile nel Cochrane database of systematic reviews) di una review prima dello studio ne limita i possibili bias

### 20.2.2 Review team

È essenziale che la revisione sia preparata da almeno due persone, per limitare possibili errori nella selezione degli studi

## 20.3 Domanda della ricerca e criteri inclusione

Per circoscrivere la ricerca bisogna aver chiaro

- PICO degli studi: nel caso non siano interventistici (es braccio singolo), comunque i gruppi da includere
  - per la popolazione: occorre definire precisamente la malattia e i suoi criteri di diagnosi, oltre eventualmente ad altri aspetti quali età, sesso, setting di cura ecc. Sui casi dubbi di studi si può evitare di escludere e decidere in un secondo momento, giustificando
  - intervento e comparazione
  - per gli outcome, vanno scelti quelli essenziali per il decision making (non più di sette) che vanno divisi in primary (non più di tre) e secondary (i rimanenti).  
Inoltre il reporting dell'outcome o meno nello studio non dovrebbe essere una causa di esclusione dello stesso (perché non fornisce dati utilizzabili, es nella metanalisi)
- tipi di studi da sintetizzare: prospettico/retrospettivo, con/senza gruppo di controllo, sperimentale o osservazionale.

*Remark 65.* Sulla tipologia di studi ci si riferirà a quanto si usa in base al quesito clinico: es per l'eziologia si procederà con sintesi di caso controllo (o coorte), diagnosi mediante studi di coorte (o trasversali o rct), terapia mediante rct, prognosi mediante coorti)

*Remark 66* (MA di studi osservazionali). Poiché le stime provenienti da studi osservazionali possono esser distorte, è verosimile che la loro combinazione sia una valutazione distorta. Ciononostante in letteratura esistono esempi di MA di studi osservazionali: i risultati sono da interpretare con cautela.

Al momento le MA di studi osservazionali non sono raccomandate in letteratura, mentre è assolutamente utile l'esecuzione di RS di studi con questo disegno; queste, pur non portando a combinazioni quantitative formali, possono fornire info sull'effetto del trattamento e possono esser utili per indirizzare la ricerca futura

## 20.4 Ricerca degli studi

Le revisioni richiedono un'approfondita e riproducibile ricerca per identificare più studi rilevanti possibili:

- coinvolgere personale esperto nella ricerca
- La ricerca solo su pubmed (MEDLINE) è generalmente considerata insufficiente
- EMBASE va considerato
- per i trial consultare il motore della cochrane (CENTRAL); sugli studi di efficacia basati su RCT ci si può forse limitare a questo
- se possibile, oltre ai primi 3, considerare altri db di letteratura: es specifici per soggetto, dottorati e dissertazioni, letteratura grigia (abstract congressi)

È importante *documentare* le ricerche effettuate; per questo e per maggiori info sulla ricerca vedere il capitolo 6 di [18].

## 20.5 Selezione degli studi e collezione dati

### 20.5.1 Selezione studi

Il tipico processo per la selezione degli studi è il seguente:

1. integrare i risultati delle ricerche (dai vari database) e rimuovere i record duplicati (stessi articoli)
2. esaminare titolo e abstract per rimuovere articoli irrilevanti (ma senza esagerare)
3. ottenere il pdf dell'articolo (per quelli rilevanti)
4. nel caso per un unico studio siano disponibili più articoli (determinarlo serve un po di lavoro da detective basato sui nomi degli autori, location e setting di studio, numero di partecipanti e dati di baseline, data/durata dello studio; mal che vada contattare gli autori) unificare il record per non considerare gli stessi risultati più volte
5. leggere gli articoli interi per determinare la compliance con i criteri di eleggibilità
6. decidere se includere o meno e procedere alla data collection sui risultati

È consigliabile che lo step 2 e a maggior ragione il 5 sia effettuato da 2 o più autori indipendentemente; consigliabile che nel team vi siano anche revisori non esperti della materia poiché gli specialisti possono avere bias sulla scelta di cosa includere o meno.

Eventuali disagreement su cosa includere o meno possono essere analizzati mediante K di Cohen (es se si classifica lo studio come includere sì/no/dubbio) e nel caso di bassi valori (poco agreement) possono rivelare la necessità di rivedere assieme i criteri di eleggibilità. Comunque il disagreement può essere solitamente risolto mediante discussione del caso o dall'arbitrato di un terzo (o infine contattando l'autore nei casi dubbi di inclusione per verificare se l'articolo rispetta i criteri)

### 20.5.2 Dati da raccogliere

L'estrazione dati sarebbe meglio farla in doppio, autonomamente per evitare errori anche qui

#### 20.5.2.1 Elementi

Vedere il manuale (capitolo 7 per approfondimenti sui singoli campi).

Source:

- Study ID (created by review author).
- Report ID (created by review author).
- Review author ID (created by review author).
- Citation and contact details.

Eligibility

- Confirm eligibility for review.
- Reason for exclusion.

Methods

- Study design.
- Total study duration.
- Sequence generation\*.
- Allocation sequence concealment\*.
- Blinding\*.
- Other concerns about bias\*.

Participants

- Total number.
- Setting.
- Diagnostic criteria.
- Age.
- Sex.
- Country.
- [Co-morbidity].
- [Socio-demographics].
- [Ethnicity].
- [Date of study].

## Interventions

- Total number of intervention groups.
- For each intervention and comparison group of interest:
- Specific intervention.
- Intervention details (sufficient for replication, if feasible).
- [Integrity of intervention].

## Outcomes

- Outcomes and time points (i) collected; (ii) reported\*.

For each outcome of interest:

- Outcome definition (with diagnostic criteria if relevant).
- Unit of measurement (if relevant).
- For scales: upper and lower limits, and whether high or low score is good.

## Results

- Number of participants allocated to each intervention group.

For each outcome of interest:

- Sample size.
- Missing participants\*.
- Summary data for each intervention group (e.g.  $2 \times 2$  table for dichotomous data; means and SDs for continuous data).
- [Estimate of effect with confidence interval; P value]z.
- [Subgroup analyses].

## Miscellaneous

- Funding source.
- Key conclusions of the study authors.
- Miscellaneous comments from the study authors.
- References to other relevant studies.
- Correspondence required.
- Miscellaneous comments by the review authors.

**20.5.2.2 Stime per misure dicotomiche**

I quattro numeri della tabella  $2 \times 2$ .



**20.5.2.3 Stime per variabili quantitative**

Media, sd e numerosità di ciascun braccio.

**20.5.2.4 Stime per analisi di sopravvivenza**

Servono le stime dei log hazard ratio e i relativi standard error. Se hanno usato cox bene, alternativamente (spesso) necessario ottenere i dati originali (a meno che non si percorrano strade più complesse, vedi 7.7.6 per dettagli)

## 20.6 Valutazione del rischio di bias negli studi inclusi

Bias è la deviazione sistematica della stima dal valore vero (non va confuso con la imprecisione delle stime); includere studi con stime affette da bias inficia la revisione.

Uno studio potrebbe essere stato eseguito coi più alti standard possibili, ma ciò nonostante non essere immune a rischi di bias (es un trial dove non è possibile effettuare il doppio cieco).

**20.6.1 Fonti di bias nei clinical trial**

L'affidabilità dei risultati di uno studio randomizzato dipendono dal fatto che le fonti di bias siano state evitate:

- *selection* bias: differenze sistematiche tra le caratteristiche di baseline dei gruppi che sono comparati; limitato dalla randomizzazione, se funziona bene
- *performance* bias: differenze sistematiche tra gruppi nella cura fornita e/o nell'esposizione a fattori oltre all'intervento di interesse; limitato da blinding dei partecipanti e del personale coinvolto nello studio
- *detection* bias: differenze tra gruppi in come l'outcome è determinato; limitato da blinding dei valutatori
- *attrition* bias: differenze tra gruppi nei ritiri dallo studio e quindi nella completezza dei dati
- *reporting* bias: differenze tra risultati riportati e non riportati (i risultati non positivi vengono più difficilmente riportati e ciò inficia i dati utilizzabili nella revisione)

## 20.7 Mantenimento della Revisione

Le revisioni sistematiche debbono essere mantenute (indicativamente ogni 2 anni) al fine di mantenere l'evidenza più aggiornata sugli effetti degli interventi. Le ragioni per un aggiornamento:

- nuovi studi disponibili

- strumenti migliori per la caratterizzazione di sottogruppi (es genetica), nuovi trattamenti, nuove misure di outcome
- nuove metodologie per la conduzione di revisioni sistematiche

## Capitolo 21

# Analisi dei dati e metanalisi

### 21.1 Outcome e misure di efficacia

*Remark 67.* Per outcome:

- dicotomici si usa OR o RR
- quantitativi si usa la differenza delle medie tra bracci o la differenza delle medie (tra bracci) standardizzata (ossia divisa la deviazione standard complessiva dell'outcome), altre volte detto *effect size*. Si ricorre al primo se tutti gli studi usano la stessa scala di misura dell'outcome, o il secondo se l'efficacia è misurata su scale diverse in studi diversi
- per tassi si usa il rate ratio
- per analisi di sopravvivenza si usa l'hazard ratio di regressioni di cox univariate

*Remark 68.* Le misure di effetto espresse come rapporto (OR, RR, rate ratio e HR) sono solitamente log trasformate (per essere simmetriche e centrate sullo 0

*Remark 69.* Allo stesso modo il display grafico di meta analisi effettuate su misure a rapporto (non precedentemente logaritmizzate) usano solitamente una scala logaritmica per affinché gli intervalli di confidenza siano simmetrici

### 21.2 Eterogeneità

*Remark 70.* I risultati di una MA sono interpretabili e utili quando gli studi che sono stati combinati erano sufficientemente comparabili, ossia *poco eterogenei*.

*Remark 71.* Vi sono due componenti principali dell'eterogeneità

**Definition 21.2.1** (Eterogeneità clinica). Si analizzano studi con differenze su pazienti, trattamento, setting di studio e outcome

**Definition 21.2.2** (Eterogeneità metodologica). Differenze su disegni (sperimentale/osservazionale) qualità e tipo di analisi (ITT/per protocol)

*Remark 72.* In generale se gli intervalli di confidenza dei singoli studi mostrano poco overlap, può essere che vi sia eterogeneità. È una condizione poco auspicabile soprattutto se gli studi sono discordi sulla direzione dell'effetto (alcuni dicono che il trattamento sperimentale sia meglio, altri che sia peggio).

*Remark 73 (Test).* Formalmente si dispone di un test che ha come ipotesi nulla quella di omogeneità degli studi: se la nulla viene rifiutata ( $p < 0.05$ ) siamo in una situazione di eterogeneità.

Dato che è un test poco potente, ogni tanto si abbassa la soglia accettando come eterogenea una situazione con  $p < 0.1$ .

Il test di eterogeneità (implementato da `revman`) è

$$Q = \sum W_i(\theta_i - \theta)^2$$

dove  $W_i$  è il peso del singolo studio,  $\theta_i$  è la stima di efficacia del singolo studio (log OR, log RR) e  $\theta$  è la stima pooled.

Sotto ipotesi nulla che non vi sia differenza negli effetti dell'intervento tra gli studi questa statistica segue una distribuzione chi quadrato con  $k - 1$  gradi di libertà e  $k$  il numero di studi analizzati

*Remark 74.* Alcuni sostengono che l'eterogeneità è inevitabile ed è dunque (non si può escluderne tutte le componenti/fonti), quindi tanto vale evitare il prendere decisioni sulla base di un test. Misure alternative che quantifichino l'eterogeneità sono state proposte

*Remark 75 ( $I^2$ ).*

$$I^2 = \left( \frac{Q - df}{Q} \right) \times 100\%$$

con  $Q$  è la statistica chi square del test di cui sopra e  $df$  sono i suoi gradi di libertà ( $k - 1$  con  $k$  studi). Questa indica la percentuale di variabilità nelle stime che è dovuta ad eterogeneità piuttosto che all'errore di campionamento (caso); una guida grezza alla sua interpretazione è

- da 0 a 40: eterogeneità potrebbe essere non importante
- da 30 a 60: potrebbe essere moderata
- da 50 a 90: potrebbe essere sostanziale
- da 70 a 100: potrebbe essere notevole

*Remark 76.* Nel caso di eterogeneità vi sono due approcci di analisi:

- analisi per sottogruppi o meta-regressioni
- stima dell'effetto del trattamento con modello ad effetti casuali (random effects model), soprattutto se quanto meno la direzione dell'effetto è abbastanza univoca

*Remark 77 (Analisi per sottogruppi).* Si può effettuare per subset di pazienti o studi

- suddivide gli studi in base ad alcune caratteristiche *pre-specificate* nel protocollo della revisione

- si effettua la meta-analisi in ogni gruppo (al fine di ridurre l'eterogeneità interna)

Per un test sulla presenza di eterogeneità si può effettuare uno test di eterogeneità sui risultati dei sottogruppi (invece che sui studi individuali).

*Remark 78* (Meta-regressione). Può quantificare l'impatto sulle stime di diverse caratteristiche dei vari studi, dovrebbe essere considerata solo se vi sono 10 studi o più

## 21.3 Metanalisi in a nutshell

*Remark 79.* Quando si fa una metanalisi non si fa una semplice somma dei pazienti e degli eventi occorsi nei singoli studi (per evitare stime inficiate dal paradosso di Simpson); si preserva invece la loro individualità e si procede ad una media pesata (al numeratore sommatoria di effetto per peso, al denominatore sommatoria dei pesi) dell'effetto dell'intervento

*Remark 80* (Fasi della metanalisi). Si procede

1. innanzitutto a calcolare una statistica per ogni studio per descrivere l'effetto del trattamento;
2. si effettua poi una stima complessiva dell'intervento calcolando una media pesata degli effetti stimati nei singoli studi

$$\frac{\sum Y_i W_i}{\sum W_i}$$

con  $Y_i$  stima dell'effetto del singolo studio (per misure rapporto come OR o RR ecc qui si usa il logaritmo) e  $W_i$  peso associato allo studio; se tutti gli studi assumono lo stesso peso, la stima è semplicemente la media degli studi

3. si può ora assumere che l'intervento abbia una sua efficacia unica/intrinseca che tutti gli studi stanno cercando di stimare e che le variazioni nelle singole stime siano dovuti all'errore di campionamento: in questo caso si procede ad una *analisi ad effetti fissi*.  
Se viceversa si pensa che l'efficacia dell'intervento abbia effettivamente una distribuzione (e non un valore singolo) si procede ad una metanalisi *ad effetti random*.
4. si calcola l'errore standard della stima complessiva, per derivare un intervallo di confidenza

## 21.4 Metodi di calcolo dei pesi $W_i$

*Remark 81* (Inverse variance method). I pesi associati a ciascuno studio sono spesso e volentieri l'inverso della varianza (o meglio del quadrato dell'errore standard) della stima per ciascuno studio

$$W_i = 1/SE^2$$

Pertanto i pesi dei vari studi saranno proporzionali alla loro dimensione: studi più grandi con errore standard tipicamente minore avranno più peso rispetto a studi più piccoli.

**Example 21.4.1.** Nella classica tabella  $2 \times 2$  con trattamento in colonna ed evento in riga

$$SE^2(\log OR) = \frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}$$

$$SE^2(\log RR) = \frac{1}{a} + \frac{1}{b} - \frac{1}{a+c} - \frac{1}{b+d}$$

**Example 21.4.2.** Nel caso di outcome quantitativo, per la differenza di medie, tipicamente

$$SE^2(meandiff) = \frac{sd_1^2}{n_1} + \frac{sd_2^2}{n_2}$$

con  $sd_1^2$  e  $n_1$  la varianza e la numerosità del primo braccio.  
Per la differenza di medie standardizzate si usa un'altra formula

**Example 21.4.3.** Per il rate ratio

$$SE^2(logratio) = \frac{1}{E_1} + \frac{1}{E_2}$$

con  $E_1$  gli eventi nel primo braccio ed  $E_2$  nel secondo

*Remark 82* (Nel caso di meta analisi con random effects). Gli errori standard studio specifici (gli SE) sono aggiustati incorporando una misura di eterogeneità tra gli effetti osservata negli studi, detta  $\tau^2$  (tau al quadrato)

*Remark 83.* In generale metodi ad effetti random ed effetti fissi daranno risultati analoghi quando non vi è eterogeneità tra gli studi.

Viceversa quando vi è eterogeneità gli intervalli di confidenza saranno più larghi nelle stime random effects rispetto a fixed effects e quindi più difficilmente si avrà un risultato overall statisticamente significativo

## Capitolo 22

# Effect size and precision

### 22.1 Overview

**treatment effects** misura di associazione tra due variabili derivante da un singolo studio che partecipa alla metanalisi in cui una rappresenta un fattore sperimentale

**effect size** misura di associazione tra due variabili di un singolo studio

**single group summary** sintesi di una variabile (quindi non una associazione di due variabili): ad esempio una prevalenza

*Remark 84.* Le metanalisi si possono analizzare qualsiasi tipo di misura, tra quelle elencate, e ai fini del metodo non cambia troppo. In generale qui si userà *effect size* in senso generico, intendendo il risultato desumibile dal singolo studio, e potendo intendere con esso anche treatment effect o single group summary

*Remark 85.* Una volta che abbiamo calcolato l'effect size e costruito mediante l'errore standard un intervallo di confidenza, le formule per calcolare un effetto complessivo, per la verifica dell'eterogeneità e così via sono le stesse indipendentemente dal tipo di effect size adoperato

### 22.2 Effect size basati su medie

Quando gli studi riportano medie e deviazioni standard l'effect size di elezione sono solitamente la differenza di medie grezze o standardizzate (o il response ratio ma mi sembra meno interessante).

#### 22.2.1 Differenza di medie non standardizzate in gruppi indipendenti

Se  $\mu_1$  e  $\mu_2$  sono le medie nelle popolazioni, la differenza in queste è  $\Delta = \mu_1 - \mu_2$ ; lo stimatore di questo parametro è

$$D = \bar{X}_1 - \bar{X}_2$$

Se non assumiamo che le due deviazioni standard del carattere nei gruppi della popolazione siano differenti, l'errore standard di  $D$  è

$$SE_D = \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

con  $S_1, S_2$  le deviazioni standard campionarie e  $n_1, n_2$  i sample size dei due gruppi

### 22.2.2 Differenza di medie standardizzate in gruppi indipendenti

Coincide con una  $d$  di Cohen e si usa tipicamente se:

- l'outcome analizzato è meno conosciuto/standard
- se studi diversi usano outcome quantitativi differenti, al fine di omogeneizzare

è definita nella popolazione come

$$\delta = \frac{\mu_1 - \mu_2}{\sigma}$$

(dove abbiamo ipotizzato che le due popolazioni siano accomunate dalla deviazione standard del carattere  $\sigma_1 = \sigma_2 = \sigma$ ).

Nel campione lo stimatore è

$$d = \frac{\bar{X}_1 - \bar{X}_2}{S_{within}}$$

con

$$S_{within} = \sqrt{\frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}}$$

con  $\bar{X}_1, \bar{X}_2$  medie campionarie,  $n_1, n_2$  i sample size e  $S_1, S_2$  le deviazioni standard nei due gruppi. Un errore standard è

$$SE_d = \sqrt{\frac{n_1 + n_2}{n_1 n_2} + \frac{d^2}{2(n_1 + n_2)}}$$

Emerge come in realtà lo stimatore  $d$  di Cohen sia lievemente biasato e tenda a sovrastimare  $\delta$  in *campioni piccoli*; il

lo stimatore corretto è il  $g$  di Hedges, che viene calcolato a partire da  $d$  e da una correzione (quella più utilizzata in pratica è la seguente)

$$J = 1 - \frac{3}{4df - 1}$$

con  $df$  sono i gradi di libertà usati per stimare  $S_{within}$  che per due gruppi indipendenti è  $n_1 + n_2 - 2$ . Si ha che lo stimatore corretto è allora

$$g = J \times d$$

e

$$SE_g = \sqrt{J^2 \times V_d}$$



	Trattati	Controlli
Eventi	A	C
Non eventi	B	D
Tot	$n_1$	$n_2$

Tabella 22.1: Tabella  $2 \times 2$ 

### 22.2.3 Response ratio

## 22.3 Effect size basati su dati binari (tabelle $2 \times 2$ )

### 22.3.1 Risk ratio

Per i risk ratio i conti (intervallo di confidenza) sono fatti su scala logaritmica (che normalizzano la distribuzione dello stimatore) e riportato sulla scala normale. Lo stimatore

$$RR = \frac{A/n_1}{C/n_2}$$

il log risk ratio è

$$LogRiskRatio = \log(RR)$$

che ha errore standard pari a

$$SE_{LogRiskRatio} = \sqrt{\frac{1}{A} - \frac{1}{n_1} + \frac{1}{C} - \frac{1}{n_2}}$$

Questo serve per costruire l'intervallo di confidenza in scala logaritmica

$$LL_{LogRiskRatio} = LogRiskRatio - 1.96 \cdot SE_{LogRiskRatio}$$

$$UL_{LogRiskRatio} = LogRiskRatio + 1.96 \cdot SE_{LogRiskRatio}$$

Per riportarlo poi il tutto in scala normale

$$LL_{RiskRatio} = \exp LL_{LogRiskRatio}$$

$$UL_{RiskRatio} = \exp UL_{LogRiskRatio}$$

### 22.3.2 Odds ratio

Anche qui si procede in scala logaritmica

$$OR = \frac{AD}{CB}$$

$$LogOddsRatio = \log(OR)$$

avente errore standard

$$SE_{LogOddsRatio} = \sqrt{\frac{1}{A} + \frac{1}{B} + \frac{1}{C} + \frac{1}{D}}$$

ed intervallo di confidenza

$$LL_{LogOddsRatio} = LogOddsRatio - 1.96 \cdot SE_{LogOddsRatio}$$

$$UL_{LogOddsRatio} = LogOddsRatio + 1.96 \cdot SE_{LogOddsRatio}$$

Per riportarlo poi il tutto in scala normale

$$LL_{OddsRatio} = \exp LL_{LogOddsRatio}$$

$$UL_{OddsRatio} = \exp UL_{LogOddsRatio}$$

### 22.3.3 Risk difference

Qui a differenza dei precedenti i calcoli su scala normale

$$RiskDiff = \frac{A}{n_1} - \frac{C}{n_2}$$

con errore standard approssimato a

$$SE_{RiskDiff} = \sqrt{\frac{AB}{n_1^3} + \frac{CD}{n_2^3}}$$

### 22.3.4 Considerazioni sulla scelta

Fare i calcoli in RR o OR poi predire la risk difference per vari risk del gruppo baseline (controllo)

## 22.4 Effect size basati su correlazioni

I conti si fanno effettuando la trasformata Z di Fisher (per avere uno stimatore approssimativamente normale); se  $\rho$  è il coefficiente di Pearson si ha la trasformata come

$$z = 0.5 \cdot \log \left( \frac{1 + \rho}{1 - \rho} \right)$$

che ha errore standard

$$SE_z = \sqrt{\frac{1}{n-3}}$$

Calcolato il CI si riporta all'unità di correlazione mediante la trasformata inversa

$$\rho = \frac{e^{2z} - 1}{e^{2z} + 1}$$

## 22.5 Conversione tra effect size

Si può fare un minestrone di misure di efficacia con dati binari, continui o di correlazione (in genere riportando in termini di  $d$  di Cohen); vedere il capitolo 7 di [6]

## Capitolo 23

# Modelli ad effetti fissi e ad effetti random

### 23.1 Introduzione

La maggior parte delle metanalisi si basa su due modelli statistici, il modello ad effetto fisso e quello ad effetti random

**Effetto fisso** si assume che vi sia un unico, vero, effect size associato al trattamento che sottosta a tutti gli studi dell'analisi considerata e che le differenze osservate tra studi siano dovuti al campionamento

**Effetti random** si ipotizza che l'effetto del trattamento (nella popolazione) possa effettivamente variare da studio a studio, quindi la variabilità che si osserva tra gli studi è dovuta sia al campionamento che al fatto che estraiamo campioni da urne differenti. Se fosse possibile effettuare infiniti studi si otterrebbe la distribuzione dell'effect size; siamo per lo più interessati ad un valore centrale (media/valore atteso) di questi effect size

*Remark 86.* La differenza principale che ne deriva riguarda la varianza, per il resto come si vedrà le formule dei due modelli sono uguali

*Remark 87.* Nella discussione che segue conviene tenere ben distinto l'*effetto vero* (quello dell'intervento nella popolazione) da quello *osservato* (ossia stimato nel campione)

*Remark 88.* In generale, di default, non vi sono motivazioni per assumere che gli effetti sottostanti (di studi differenti, in setting lievemente differenti ecc) siano coincidenti; pertanto, di solito, random effects all the way!

L'analisi ad effetti fissi ha senso se due condizioni sono rispettate: tutti gli studi della metanalisi sono funzionalmente identici e il nostro obiettivo è stimare l'effect size comune per questo tipo di popolazione e non si cerca di generalizzare ad altre popolazioni. Se viceversa i ricercatori stanno accumulando dati da studi che sono stati eseguiti da ricercatori che hanno operato indipendentemente, sarebbe inverosimile che tutti gli studi fossero funzionalmente equivalenti.

Alcuni effettuano una analisi ad effetti random dopo avere effettuato un test di eterogeneità ma non è la strada che si consiglia qui; viceversa se uno fa

l'analisi a effetti fissi poi l'eterogeneità ne esce, allora l'ipotesi su cui si è basato forse non regge.

## 23.2 Effetto fisso

Sotto effetto fisso ipotizziamo che vi sia un unico effect size associato al trattamento,  $\theta$  e che ogni effect size rilevato in ciascuno studio incluso nell'analisi  $Y_i$  differisca di un errore casuale  $\varepsilon_i$  dovuto al campione estratto

$$Y_i = \theta + \varepsilon_i \quad (23.1)$$

In una analisi ad effetti fissi ad ogni studio viene assegnato un peso pari a

$$W_i = \frac{1}{V_{Y_i}} \quad (23.2)$$

con  $V_{Y_i}$  la cd varianza (il quadrato dell'errore standard dello stimatore considerato applicato allo studio). La stima dell'effetto complessivo per una metanalisi con  $k$  studi è dato da

$$M = \frac{\sum_{i=1}^k W_i Y_i}{\sum_{i=1}^k W_i} \quad (23.3)$$

La varianza di questa stima complessiva è data da

$$V_M = \frac{1}{\sum_{i=1}^k W_i} \quad (23.4)$$

quindi è di fatto legata direttamente alla varianza dei singoli studi. L'errore standard dello stimatore complessivo è la radice di questa

$$SE_M = \sqrt{V_M} \quad (23.5)$$

Ottenute queste stime si può costruire l'intervallo di confidenza dello stimatore complessivo

$$LL = M - 1.96 \cdot SE_M \quad UL = M + 1.96 \cdot SE_M \quad (23.6)$$

o procedere a test Z sulla presenza di effetto come abitualmente mediante

$$Z = \frac{M}{SE_M} \quad (23.7)$$

e confrontando con i quantili soglia della normale standard.

## 23.3 Effetti random

Come detto si ipotizza che vi sia una distribuzione di effect size piuttosto che un singolo effect size unico.

In una analisi ad effetti random si ipotizza che l'effect size di un singolo studio si discosti dal valore medio/atteso degli effect size per questo tipo di intervento per due componenti:

- innanzitutto l'effect size della popolazione/trattamento considerata si discosta dal valore atteso degli effect size per una componente  $\zeta_i$  (se avessimo un numero di studi infiniti potremmo calcolare l'effect size atteso complessivo)
- l'effect size osservato nel campione si discosta dall'effect size della popolazione considerata sempre per un errore casuale (se avessimo un campione infinito l'errore sarebbe nullo)

In altre parole si ha che l'effect size osservato sia la somma di tre componenti:

$$Y_i = \underbrace{\mu + \zeta_i}_{\theta_i} + \varepsilon_i \quad (23.8)$$

con  $\mu$  il valore atteso degli effect size complessivamente e  $\theta_i$  l'effect size della popolazione considerata.

La distanza di  $\theta_i$  da  $\mu$  dipende dalla variabilità della distribuzione degli effetti tra gli studi, indicata con  $\tau$  (se si pensa alla deviazione standard) o  $\tau^2$  (se si pensa alla varianza. Essendo un valore unico per tutto la metanalisi, rientra nei conteggi della variabilità di tutti gli studi considerati, come si vedrà.

Uno stimatore di  $\tau^2$  (Dersimonian e Laird) è il seguente:

$$T^2 = \frac{Q - df}{C} \quad (23.9)$$

dove

$$\begin{aligned} Q &= \sum_{i=1}^k W_i Y_i^2 - \frac{\left(\sum_{i=1}^k W_i Y_i\right)^2}{\sum_{i=1}^k W_i} \\ df &= k - 1 \\ C &= \sum_i W_i - \frac{\sum W_i^2}{\sum W_i} \end{aligned}$$

con  $k$  il numero di studi.

Una volta ottenuta la stima, la stima della varianza per il singolo studio diviene

$$V_{Y_i}^* = V_{Y_i} + T^2$$

ossia si aggiunge questa componente a quella già adottata per il modello a effetto fisso. I calcoli poi procedono analogamente in tutto e per tutto a quanto già visto in tal caso.

Si approfondirà il discorso quando di parlerà di eterogeneità

## 23.4 Un confronto

*Remark 89.* A parità di altre condizioni, i modelli ad effetti per la presenza di una componente comune di varianza (positiva, che viene sommata ad un'altra componente positiva) rende complessivamente le varianze dei vari studi più simili tra loro, dunque i pesi assegnati più omogenei e complessivamente rispetto ad un'analisi ad effetto fisso pesa maggiormente gli studi più piccoli e meno quelli

grossi; a parole dato che ogni studio fornisce informazioni su differenti effect size, vogliamo assicurarci che tutti siano rappresentati nella stima complessiva e che non vogliamo penalizzare alcuni effect size solo perché i relativi studi siano piccoli (come si fa in uno studio ad effetto fisso); per lo stesso ragionamento non diamo eccessivo peso a studi con molti pazienti.

*Remark 90.* In merito all'interpretazione dell'intervallo di confidenza, nel caso dell'effetto fisso si tratta dell'intervallo di confidenza di questo, mentre in quello di effetti random si tratta dell'intervallo di confidenza del valore atteso dei vari effect size

*Remark 91.* Sull'interpretazione del test, per l'effetto fisso l'ipotesi nulla è che l'effetto sia nullo, mentre in quello degli effetti random è che il valore atteso/medio degli effetti random sia nullo (ma non si esclude che vi possano essere popolazioni/setting in cui l'effect size sia effettivamente diverso)

## 23.5 Esempi

Qui si riproducono gli esempi del capitolo 14 di [6] mediante il pacchetto `metafor` di R. Dato che la differenza principale si ha solamente nel calcolo degli effect size si approfondisce per esteso il caso di dati dicotomici, facendo vedere quello che varia nella stima per i continui e correlazionali.

### 23.5.1 Dati dicotomici

Procediamo alla stima degli effect size di ciascuno studio mediante `escalc` e poi alla metanalisi mediante `rma` (o `rma.uni`, sinonimo)

```
library(lbdatasets)
library(metafor)
metabin

##      study te tne  tn ce cne  cn
## 1  Saint 12  53  65 16  49  65
## 2  Kelly  8  32  40 10  30  40
## 3 Pilbeam 14  66  80 19  61  80
## 4   Lane 25 375 400 80 320 400
## 5 Wright  8  32  40 11  29  40
## 6   Day 16  49  65 18  47  65

## questo fornisce il calcolo degli effect size per
## il caso degli odds-ratio (stime sulla log scale)
## cfr pag93 tab 14.5
(es_bin <- escalc(measure = 'OR', # log odds-ratio
                 ai = te,        # eventi nei trattati
                 bi = tne,       # non eventi nei trattati
                 n1i = tn,       # per check totale trattati
                 ci = ce,        # eventi nei controlli
                 di = cne,       # non eventi nei controlli
                 n2i = cn,       # per check totale controlli
                 data = metabin)) # data.frame di riferimento
```

```
##
##      study te tne  tn ce cne  cn      yi      vi
## 1   Saint 12  53  65 16  49  65 -0.3662 0.1851
## 2   Kelly  8  32  40 10  30  40 -0.2877 0.2896
## 3 Pilbeam 14  66  80 19  61  80 -0.3842 0.1556
## 4    Lane 25 375 400 80 320 400 -1.3218 0.0583
## 5 Wright  8  32  40 11  29  40 -0.4169 0.2816
## 6    Day 16  49  65 18  47  65 -0.1595 0.1597

## fixed effect meta analysis
(bin_fe <- rma(yi = yi, vi = vi, data = es_bin, method = 'FE'))

##
## Fixed-Effects Model (k = 6)
##
## I^2 (total heterogeneity / total variability):  52.61%
## H^2 (total variability / sampling variability):  2.11
##
## Test for Heterogeneity:
## Q(df = 5) = 10.5512, p-val = 0.0610
##
## Model Results:
##
## estimate      se      zval      pval      ci.lb      ci.ub
## -0.7241  0.1539  -4.7068  <.0001  -1.0257  -0.4226  ***
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

## random effects meta analysis (Tau^2 stimato secondo dersimonian e laird)
(bin_re <- rma(yi = yi, vi = vi, data = es_bin, method = 'DL'))

##
## Random-Effects Model (k = 6; tau^2 estimator: DL)
##
## tau^2 (estimated amount of total heterogeneity): 0.1729 (SE = 0.2148)
## tau (square root of estimated tau^2 value):      0.4158
## I^2 (total heterogeneity / total variability):  52.61%
## H^2 (total variability / sampling variability):  2.11
##
## Test for Heterogeneity:
## Q(df = 5) = 10.5512, p-val = 0.0610
##
## Model Results:
##
## estimate      se      zval      pval      ci.lb      ci.ub
## -0.5663  0.2388  -2.3711  0.0177  -1.0344  -0.0982  *
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Una volta ottenuta la stima di metanalisi possiamo indagare l'oggetto; concentrandoci sul caso degli effetti random ...

```
## cosa c'è
names(bin_re)

## [1] "b"          "beta"       "se"         "zval"       "pval"
## [6] "ci.lb"      "ci.ub"      "vb"         "tau2"       "se.tau2"
## [11] "tau2.fix"   "tau2.f"     "I2"         "H2"         "R2"
## [16] "vt"         "QE"         "QEp"        "QM"         "QMdf"
## [21] "QMp"        "k"          "k.f"        "k.eff"      "k.all"
## [26] "p"          "p.eff"      "parms"      "int.only"   "int.incl"
## [31] "intercept"  "allvupos"   "coef.na"    "yi"         "vi"
## [36] "X"          "weights"    "yi.f"       "vi.f"       "X.f"
## [41] "weights.f"  "M"          "chksumyi"   "chksumvi"   "chksumX"
## [46] "outdat.f"   "ni"         "ni.f"       "ids"        "not.na"
## [51] "subset"     "slab"       "slab.null"  "measure"    "method"
## [56] "model"      "weighted"   "test"       "dfs"        "ddf"
## [61] "s2w"        "btt"        "m"          "digits"     "level"
## [66] "control"    "verbose"    "add"        "to"         "drop00"
## [71] "fit.stats"  "data"       "formula.yi" "formula.mods" "version"
## [76] "call"       "time"

## alcune stime puntuali presentate tra pag 96 e 97
with(bin_re, c('stima' = b,
               'se' = se,
               'ci.low' = ci.lb,
               'ci.up' = ci.ub,
               'z' = zval,
               'two_tail_p' = pval))

##          stima          se      ci.low      ci.up          z      two_tail_p
## -0.56629590  0.23883443 -1.03440279 -0.09818902 -2.37108149  0.01773612

## per ottenere i pesi relativi (sommano a 100)
weights(bin_re)

##          1          2          3          4          5          6
## 15.93285 12.33370 17.36382 24.67247 12.54918 17.14797

## stime e ci di alcuni pararametri stimati
confint(bin_re)

##
##          estimate  ci.lb  ci.ub
## tau^2      0.1729 0.0000 0.9656
## tau        0.4158 0.0000 0.9826
## I^2(%)     52.6118 0.0000 86.1112
## H^2        2.1102 1.0000 7.2001
```

Infine procediamo alla visualizzazione del forestplot (sulla scala degli odds-



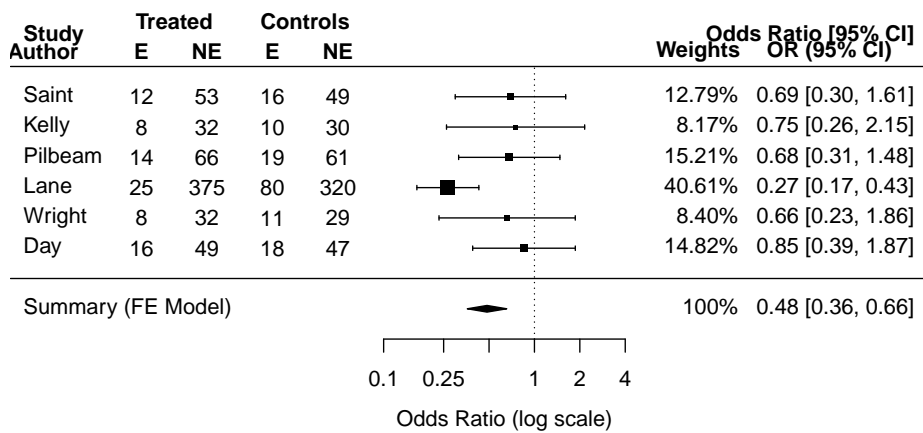


Figura 23.1: Forestplot metanalisi fixed effect

ratio e cercando di ottimizzare un minimo) dell'analisi ad effetto fisso in figura 23.1 mentre quella con effetti random in figura 23.2.

```
## -----
## Fixed effect MA forest plot
## -----
par(mar = c(5,1,1,1))
data_matrix <- with(metabin, cbind(te, tne, ce, cne))
data_matrix_col_pos <- seq(-6, -3)
ticks_pos <- log(c(0.1, 0.25, 0.5, 1, 2, 4))
xlim <- c(-8, 6)
forest(x = bin_fe,      # fornire la stima degli effect size
      atransf = exp,    # per ottenere gli OR invece dei log (OR)
      at = ticks_pos,  # posizionamento dei ticks rispettando la scala
      xlim = xlim,
      ## matrice dei dati a sinistra e posizionamento delle colonne
      ilab = data_matrix,
      ilab.xpos = data_matrix_col_pos,
      slab = metabin$study, ## etichette di riga
      showweights = TRUE,  ## display dei pesi
      mlab = 'Summary (FE Model)')
## intestazioni
par(font = 2)
title_first_row <- 8.5
title_second_row <- 7.5
text(x = data_matrix_col_pos, y = title_second_row,
     labels = rep(c("E", "NE"), 2))
text(x = -7.5, y = title_second_row, 'Author')
text(x = c(-5.5, -3.5), y = title_first_row, c('Treated', 'Controls'))
text(x = c(2.5, 4.5), y = title_second_row, c('Weights', 'OR (95% CI)'))
```

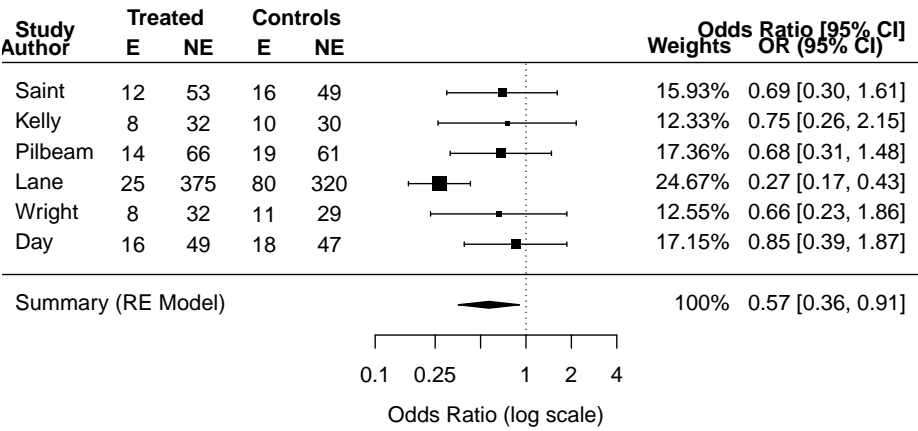


Figura 23.2: Forestplot metanalisi random effects

```
## -----
## Random effects MA forest plot
## -----
par(mar = c(5,1,1,1))
forest(x = bin_re,      # fornire la stima degli effect size
      atransf = exp,    # per ottenere gli OR invece dei log (OR)
      at = ticks_pos,  # posizionamento dei ticks rispettando la scala
      xlim = xlim,
      ## matrice dei dati a sinistra e posizionamento delle colonne
      ilab = data_matrix,
      ilab.xpos = data_matrix_col_pos,
      slab = metabin$study, ## etichette di riga
      showweights = TRUE,  ## display dei pesi
      mlab = 'Summary (RE Model)')
## intestazioni
par(font = 2)
text(x = data_matrix_col_pos, y = title_second_row,
     labels = rep(c("E", "NE"), 2))
text(x = -7.5, y = title_second_row, 'Author')
text(x = c(-5.5, -3.5), y = title_first_row, c('Treated', 'Controls'))
text(x = c(2.5, 4.5), y = title_second_row, c('Weights', 'OR (95% CI)'))
```

23.5.2 Dati continui

```
metacon

##      study tm tsd  tn cm csd  cn
## 1 Carroll 94  22  60 92  20  60
## 2  Grant 98  21  65 92  22  65
## 3   Peck 98  28  40 88  26  40
```

```
## 4 Donat 94 19 200 82 17 200
## 5 Stewart 98 21 50 88 22 45
## 6 Young 96 21 85 92 22 85

## cfr tab 14.2 pag 88
(es_con <- escalc(measure = 'SMD', # SMD
                 m1i = tm, # media trattati
                 sd1i = tsd, # sd trattati
                 n1i = tn, # n trattati
                 m2i = cm, # media controllati
                 sd2i = csd, # sd controllati
                 n2i = cn, # n controllati
                 data = metacon)) # data.frame di riferimento

##
## study tm tsd tn cm csd cn yi vi
## 1 Carroll 94 22 60 92 20 60 0.0945 0.0334
## 2 Grant 98 21 65 92 22 65 0.2774 0.0311
## 3 Peck 98 28 40 88 26 40 0.3665 0.0508
## 4 Donat 94 19 200 82 17 200 0.6644 0.0106
## 5 Stewart 98 21 50 88 22 45 0.4618 0.0433
## 6 Young 96 21 85 92 22 85 0.1852 0.0236

## ## fixed effect meta analysis
## (con_fe <- rma(yi = yi, vi = vi, data = es_con, method = 'FE'))
## ## random effects meta analysis (Tau^2 stimato secondo dersimonian e laird)
## (con_re <- rma(yi = yi, vi = vi, data = es_con, method = 'DL'))
```

### 23.5.3 Correlazioni

```
metacor

## study cor n
## 1 Fonda 0.50 40
## 2 Newman 0.60 90
## 3 Grant 0.40 25
## 4 Granger 0.20 400
## 5 Milland 0.70 60
## 6 Finch 0.45 50

## vedere per confronto table 14.8 pag 98
(es_cor <- escalc(measure = 'ZCOR', # trasformata di Fisher su correlazione
                 ri = cor, # correlazioni raw
                 ni = n, # n per gruppo
                 data = metacor)) # data.frame di riferimento

##
## study cor n yi vi
## 1 Fonda 0.50 40 0.5493 0.0270
```

```
## 2  Newman 0.60  90 0.6931 0.0115
## 3   Grant 0.40  25 0.4236 0.0455
## 4 Granger 0.20 400 0.2027 0.0025
## 5 Milland 0.70  60 0.8673 0.0175
## 6   Finch 0.45  50 0.4847 0.0213

## ## fixed effect meta analysis
## (cor_fe <- rma(yi = yi, vi = vi, data = es_cor, method = 'FE'))
## ## random effects meta analysis (Tau^2 stimato secondo dersimonian e laird)
## (cor_re <- rma(yi = yi, vi = vi, data = es_cor, method = 'DL'))
```

## Capitolo 24

# Eterogeneità

### 24.1 Quantificazione

Per eterogeneità intendiamo la variabilità della distribuzione dell'effect size nella popolazione: sotto un modello fixed essa si assume essere 0 mentre sotto un modello random un valore positivo.

Nel caso vi sia eterogeneità, la variabilità complessiva degli effect size può/deve essere spaccettata in due componenti:

- variabilità dovuta all'eterogeneità
- variabilità dovuta al campionamento casuale

Un primo indice di eterogeneità di  $k$  studi è

$$Q = \sum_{i=1}^k W_i (Y_i - M)^2 = \sum_{i=1}^k \left( \frac{Y_i - M}{S_i} \right)^2 \quad (24.1)$$

con

- $W_i = 1/V_i$  peso dello studio e  $V_i$  sua varianza (e  $S_i$  errore standard)
- $Y_i$  l'effect size del singolo studio
- $M$  l'effect size complessivo derivante dai  $k$  studi

Si tratta di una somma di scarti, standardizzati al quadrato e misura la variabilità totale presente.

#### 24.1.1 Test di eterogeneità

Se vogliamo un test di eterogeneità, nell'ipotesi nulla di assenza di eterogeneità (es se  $Y_i = M + \varepsilon_i$  con  $\varepsilon_i$  normale), il  $Q$  si distribuisce come un  $\chi^2$  con  $df = k - 1$  gradi di libertà. Anche qui, se il test è non significativo non vuol necessariamente dire che non vi sia eterogeneità, potrebbe essere low power (per questo spesso si usa invece di 0.05 0.1 come soglia decisionale)

### 24.1.2 Scarto di eterogeneità

Possiamo poi produrre tutto un altro range di indicatori partendo dalla differenza tra il valore osservato di variabilità totale  $Q$  e quello atteso in assenza di eterogeneità.

Nel caso in cui tutti gli studi condividano un effect size comune, il valore atteso di  $Q$  è una somma di scarti (determinati dall'errore di campionamento casuale) pari ai gradi di libertà

$$df = k - 1$$

e la misura

$$Q - df$$

se positiva indicherà l'eccesso di variazione dovuto alla eterogeneità, mentre se negativa indicherà una scarsissima eterogeneità (minore di quella ragionevolmente assumibile come casuale).

### 24.1.3 Stima di $\tau^2$

Con  $\tau^2$  si intende la varianza degli effetti associati al trattamento nella popolazione; dato che non la possiamo stimare direttamente, ricorriamo ad uno stimatore  $T^2$  basato sugli studi a nostra disposizione che si ottiene come

$$T^2 = \frac{Q - df}{C} \quad (24.2)$$

con

$$C = \sum_i W_i - \frac{\sum W_i^2}{W_i} \quad (24.3)$$

$T^2$  è usata per il calcolo dei pesi nei modelli random effect come

$$W_i^* = \frac{1}{V_{Y_i}^*} = \frac{1}{V_{Y_i} + T^2} \quad (24.4)$$

### 24.1.4 $I^2$

È un indicatore della percentuale di variabilità complessiva dovuta ad eterogeneità e calcola come

$$I^2 = \frac{Q - df}{Q} \times 100 \quad (24.5)$$

dove al numeratore abbiamo un indicatore di  $\tau^2$  al denominatore un indicatore della varianza totale. Vi sono anche intervalli di confidenza per  $I^2$ , volendo. Higgins suggerisce (pag 119 libro) di considerare 25, 50 e 75 come eterogeneità bassa, media ed elevata

### 24.1.5 Applicazioni in R

In metafor queste misure vengono calcolate (alcune soltanto nel caso di modelli modelli random effects) mediante `rma`, come si è visto negli esempi precedenti

## 24.2 Prediction intervals

Per le stime globali associate ad una metanalisi possiamo provvedere un intervallo di confidenza del valore centrale, che ha come scopo quello di descrivere un set di valori verosimili per l'effetto complessivo (stima fixed) o dell'effetto medio (stima random).

Nel caso di modelli random, ove abbiamo una stima di  $\tau^2$ , possiamo anche fornire un intervallo di predizione ossia (specularmente ad un intervallo di predizione su una retta di regressione di uno studio primario) un range di valori verosimili per il valore di effetto di un nuovo studio. Si costruisce come

$$\mu \pm z_{1-\alpha/2} \sqrt{\tau^2}$$

ma se non si hanno i valori della popolazione qui richiesti si rimedia con

$$M^* \pm t_{1-\alpha/2, df} \sqrt{T^2 + V_{M^*}}$$

dove  $M^*$  è l'effect size medio e  $V_{M^*}$  la sua varianza.

Graficamente l'intervallo di predizione si espande come una barretta oltre il diamante dell'intervallo di confidenza della media

## 24.3 Analisi per sottogruppi

Capitolo 19 libro

TODO: fixme

## 24.4 Metaregressione

La regressione può esser utilizzata anche avendo come unità di analisi il singolo studio (nello specifico il suo effect size, spesso e volentieri logaritmizzato per renderlo simmetrico attorno allo zero). Valgono anche qui considerazioni sul rapporto tra covariate adottate e studi/pazienti disponibili (anche se non vi sono regole scritte nella pietra).

Si studia l'esempio del vaccino BCG nel prevenire la tubercolosi, analisi fatta sui risk ratio (si riproduce solamente la stima con random effects)

```
options(width = 100)
library(metafor)
(db <- dat.bcg) # dataset impiegato, disponibile in metafor
```

##	trial	author	year	tpos	tneg	cpos	cneg	ablat	alloc
## 1	1	Aronson	1948	4	119	11	128	44	random
## 2	2	Ferguson & Simes	1949	6	300	29	274	55	random
## 3	3	Rosenthal et al	1960	3	228	11	209	42	random
## 4	4	Hart & Sutherland	1977	62	13536	248	12619	52	random
## 5	5	Frimodt-Moller et al	1973	33	5036	47	5761	13	alternate
## 6	6	Stein & Aronson	1953	180	1361	372	1079	44	alternate
## 7	7	Vandiviere et al	1973	8	2537	10	619	19	random
## 8	8	TPT Madras	1980	505	87886	499	87892	13	random
## 9	9	Coetzee & Berjak	1968	29	7470	45	7232	27	random
## 10	10	Rosenthal et al	1961	17	1699	65	1600	42	systematic

```
## 11      11      Comstock et al 1974 186 50448 141 27197 18 systematic
## 12      12      Comstock & Webster 1969 5 2493 3 2338 33 systematic
## 13      13      Comstock et al 1976 27 16886 29 17825 33 systematic

rr <- escalc(measure = 'RR', # log
             ai = tpos,
             bi = tneg,
             ci = cpos,
             di = cneg,
             data = db)
rr <- rr[with(rr, order(yi)), ]
rr[, names(rr) %without% c("alloc")]

##
##      trial      author year tpos  tneg cpos  cneg ablat      yi      vi
## 7          7      Vandiviere et al 1973 8 2537 10 619 19 -1.6209 0.2230
## 2          2      Ferguson & Simes 1949 6 300 29 274 55 -1.5854 0.1946
## 4          4      Hart & Sutherland 1977 62 13536 248 12619 52 -1.4416 0.0200
## 10         10      Rosenthal et al 1961 17 1699 65 1600 42 -1.3713 0.0730
## 3          3      Rosenthal et al 1960 3 228 11 209 42 -1.3481 0.4154
## 1          1      Aronson 1948 4 119 11 128 44 -0.8893 0.3256
## 6          6      Stein & Aronson 1953 180 1361 372 1079 44 -0.7861 0.0069
## 9          9      Coetzee & Berjak 1968 29 7470 45 7232 27 -0.4694 0.0564
## 11         11      Comstock et al 1974 186 50448 141 27197 18 -0.3394 0.0124
## 5          5      Frimodt-Moller et al 1973 33 5036 47 5761 13 -0.2175 0.0512
## 13         13      Comstock et al 1976 27 16886 29 17825 33 -0.0173 0.0714
## 8          8      TPT Madras 1980 505 87886 499 87892 13 0.0120 0.0040
## 12         12      Comstock & Webster 1969 5 2493 3 2338 33 0.4459 0.5325

## Stima overall FE, compliant col libro
(ma <- rma(yi = yi, vi = vi, method = 'DL', data = rr))

##
## Random-Effects Model (k = 13; tau^2 estimator: DL)
##
## tau^2 (estimated amount of total heterogeneity): 0.3088 (SE = 0.2299)
## tau (square root of estimated tau^2 value): 0.5557
## I^2 (total heterogeneity / total variability): 92.12%
## H^2 (total variability / sampling variability): 12.69
##
## Test for Heterogeneity:
## Q(df = 12) = 152.2330, p-val < .0001
##
## Model Results:
##
## estimate      se      zval      pval      ci.lb      ci.ub
## -0.7141 0.1787 -3.9952 <.0001 -1.0644 -0.3638 ***
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



```
## Metaregressione FE: si pone il rhs della formula in mods (ablat = latitude)
(mr <- rma(yi = yi, vi = vi, mods = ~ ablat, method = 'DL', data = rr))

##
## Mixed-Effects Model (k = 13; tau^2 estimator: DL)
##
## tau^2 (estimated amount of residual heterogeneity):      0.0633 (SE = 0.0548)
## tau (square root of estimated tau^2 value):             0.2516
## I^2 (residual heterogeneity / unaccounted variability): 64.21%
## H^2 (unaccounted variability / sampling variability):    2.79
## R^2 (amount of heterogeneity accounted for):             79.50%
##
## Test for Residual Heterogeneity:
## QE(df = 11) = 30.7331, p-val = 0.0012
##
## Test of Moderators (coefficient 2):
## QM(df = 1) = 18.8452, p-val < .0001
##
## Model Results:
##
##           estimate      se      zval      pval      ci.lb      ci.ub
## intrcpt      0.2595  0.2323   1.1172  0.2639   -0.1958    0.7149
## ablat       -0.0292  0.0067  -4.3411 <.0001   -0.0424   -0.0160 ***
##
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



Parte VIII

Dimensionamento  
campionario



## Capitolo 25

# Introduzione al dimensionamento campionario

### 25.1 Approcci e ambiti di dimensionamento

#### 25.1.1 Errori nei test di ipotesi

Distinguiamo:

- errore di primo tipo: indicato con  $\alpha$  è la probabilità di rifiutare l'ipotesi nulla quando essa in realtà è vera (trattamento non funziona ma concludiamo erroneamente che è efficace);
- errore del secondo tipo: indicato con  $\beta$  è la probabilità di non rifiutare la nulla quando essa è in realtà falsa (il trattamento è efficace ma concludiamo erroneamente il contrario)

#### 25.1.2 Giustificazione del dimensionamento

L'attività di dimensionare propriamente uno studio è volta ad evitare:

- studi sottodimensionati: con poche probabilità di dimostrare una differenza, qualora sia presente, sono di fatto poco etici poiché si sottopongono pazienti ad uno studio senza che ve ne sia un beneficio scientifico;
- studi sovradimensionati: oltre ad essere costosi sono poco etici poiché se con meno pazienti potessimo dimostrare l'efficacia del trattamento non esporremmo i pazienti in eccesso del gruppo di controllo ad un trattamento subottimale.

#### 25.1.3 Approcci al dimensionamento

Due approcci possibili, complementari e che rispondono a necessità differenti, al dimensionamento[11]:

- *precision analysis*: si vuole determinare quanti pazienti sono necessari al fine di ottenere stime con la precisione (es ampiezza intervallo di confidenza)

reputata ragionevole/di interesse (maggior enfasi sul controllo dell'errore di primo tipo);

- *power analysis* si vuole determinare quanti soggetti sono necessari (in un campione) per avere una determinata probabilità (es 80%) di identificare (ossia test statisticamente significativi) una data differenza, reputata clinicamente rilevante, qualora essa esista effettivamente (nella popolazione). (maggior enfasi sul controllo dell'errore del secondo tipo)

### 25.1.4 Ambiti di dimensionamento

Dimensionamento può essere fatto in 4 ambiti differenti:

- *stima*: caso classico dove si calcolano i soggetti necessari in base alla domanda
- *giustificazione*: dove si fornisce una giustificazione per un campione selezionato a priori (es sulla base di fattibilità)
- *adjustment*: dove il sample size stimato in precedenza viene aggiustato per fattori come dropout o covariate, al fine di fornire un numero sufficiente per l'analisi da condurre
- *re-estimation*: nell'ambito degli studi con analisi ad interim si provvede ad aggiustamento basato sulle informazioni raccolte

## 25.2 Ipotesi a confronto e disegni

Per fornire un calcolo affidabile deve essere scelto in anticipo un test per l'ipotesi di interesse, determinato in base al disegno dello studio e alla domanda alla quale vuole rispondere. Spesso il focus principale è l'efficacia o la *safety*, e per entrambe si possono configurare confronti di [11]:

- *eguaglianza*: le ipotesi al confronto sono:

$$H_0 : \mu_T = \mu_C \quad \text{vs} \quad H_1 : \mu_T \neq \mu_C \quad (25.1)$$

con  $\mu_C, \mu_T$  la risposta nella variabile outcome per controlli (o valore teorico) e trattati rispettivamente. Il rifiuto della nulla suggerisce vi sia differenza tra trattati e controlli;

- *non inferiorità di un margine*: le ipotesi al confronto sono

$$H_0 : \mu_C - \mu_T \geq \delta \quad \text{vs} \quad H_1 : \mu_C - \mu_T < \delta \quad (25.2)$$

dove  $\delta$  è una differenza di interesse clinico ritenuta importante. Il rifiuto della nulla suggerisce che la differenza tra trattati e controlli sia inferiore ad una differenza rilevante  $\delta$ , e quindi il trattamento sperimentale sia efficace quanto la terapia standard (comune in trial di efficacia dove il trattamento sperimentale è meno tossico dello standard, più facile da amministrare o meno costoso);

- *superiorità di un margine*: le ipotesi al confronto sono

$$H_0 : \mu_T - \mu_C \leq \delta \quad \text{vs} \quad H_1 : \mu_T - \mu_C > \delta \quad (25.3)$$

Nel caso si rifiuti la nulla la differenza tra trattati e controlli è maggiore di una soglia reputata rilevante e in questo senso il trattamento sperimentale è superiore rispetto alla terapia standard.

Da notare che le ipotesi di cui sopra testano la cosiddetta *superiorità clinica*; nel caso  $\delta = 0$  alle ipotesi di sopra ci si riferisce come di *superiorità statistica*

- *equivalenza*; le ipotesi confrontate sono

$$H_0 : |\mu_T - \mu_C| \geq \delta \quad \text{vs} \quad H_1 : |\mu_T - \mu_C| < \delta \quad (25.4)$$

Nel caso la nulla venga rifiutata, si conclude che la differenza tra trattati e controlli non sia clinicamente rilevante

Le ipotesi di confronto (eguaglianza, ecc) debbono essere ben chiare all'atto del dimensionamento poiché il dimensionamento (formule) dipende da esse.

## 25.3 Considerazioni assortite

### 25.3.1 Test a una o due code

Nell'ambito della power analysis lo studio può prevedere un test ad una o due code:

- il beneficio del test ad una coda, a parità di altre condizioni è che richiede tipicamente un campione inferiore
- lo svantaggio è che ci si espone maggiormente ad un rischio di falsi positivi; nel caso non vi sia efficacia del trattamento commettiamo un errore del primo tipo nel 5% dei casi con un test ad una coda, solo nel 2.5% dei casi per un test a due code. Nella realtà se per l'approvazione di un trattamento occorrono 2 trial indipendenti, il rischio di falsi positivi nel caso questi trial usino test a una coda è  $0.05^2 = 0.0025 = 0.25\%$ , mentre nel caso di due trial che impieghino test a due code il rischio di falsi positivi scende a  $0.025^2 = 0.000625 = 0.0625\%$ .

Alcuni ricercatori reputano 0.25% un rischio comunque accettabile, giustificando l'impiego di test a una coda, l'FDA sembra preferire test a due code

### 25.3.2 Aggiustamenti per dropouts

In presenza di dropout e dati per l'outcome principale

**Example 25.3.1.** Se il calcolo del sample size ci suggerisce 86 pazienti, ma si reputa verosimile un dropout del 20% sarà necessario arruolare

$$(86/(100 - 20)) \cdot 100 = 86/0.8 = 107.5 \approx 108$$

pazienti

### 25.3.3 Pacchetti R

Pacchetti utili:

- **TrialSize** implementa le funzioni per [11]
- **presize** calcola il campione sulla base di stima e ampiezza dell'intervallo di confidenza (oppure l'ampiezza garantita da un campione)
- **CRTSize** sample size per cluster randomized trials
- **clinfun** ha funzioni per dimensionamento e analisi di studi di fase 2, test esatto di Fisher
- **CRM** ha un Continual Reassessment Method per le fasi 1



# Capitolo 26

## Un gruppo

### 26.1 Precision analysis - casi base

In questi casi si determina il campione basandosi sull'errore di primo tipo e utilizzando gli approcci degli intervalli di confidenza.

La precisione di una stima dipende dall'ampiezza del suo intervallo di confidenza:

1. direttamente dal livello di confidenza ad  $(1 - \alpha) \cdot 100\%$
2. direttamente dalla variabilità del fenomeno
3. inversamente dal numero di soggetti impiegati nella stima

Essendo tipicamente i primi due parametri considerati fissi/dati si agisce sul terzo al fine di avere una precisione di stima (ampiezza dell'intervallo di confidenza) accettabile.

**Definition 26.1.1** (Errore massimo (Maximum error)). Si chiama così la semiampiezza (ampiezza/2) massima dell'intervallo di confidenza che si è disposti ad accettare.

#### 26.1.1 Stima di una media

Nel caso  $n$  iid normali  $y_1, \dots, y_n$  con media  $\mu$  e varianza  $\sigma^2$ .

##### 26.1.1.1 Intervallo a due code

**Varianza della popolazione nota** Qualora la varianza sia nota un intervallo di confidenza a due code per  $\mu$ , con un livello di confidenza pari a  $(1 - \alpha) \cdot 100\%$  è dato da:

$$\hat{\mu} \pm z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \quad (26.1)$$

L'errore massimo che siamo disposti a commettere è

$$E = |\hat{\mu} - \mu| = z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} \quad (26.2)$$

e risolvendo per  $n$  si ottiene l'ampiezza campionaria in grado di garantirlo

$$n = \frac{z_{1-\alpha/2}^2 \cdot \sigma^2}{E^2} \quad (26.3)$$

Si noti come in questo approccio non si fa uso dell'errore  $\beta$ .

**TODO:** dire meglio qui

**Example 26.1.1.** Vogliamo determinare il campione necessario per avere il 95% di probabilità che l'errore nella stima effettuata sia meno del 10% della deviazione standard del fenomeno (ossia  $0.1\sigma$ ). Si ha che l'errore massimo è

$$z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}} = 0.1\sigma$$

E pertanto

$$n = \frac{z_{1-\alpha/2}^2 \sigma^2}{E^2} = \frac{(1.96)^2 \sigma^2}{(0.1\sigma)^2} = 384.2 \approx 385$$

**Varianza ignota** Nel caso in cui la varianza non sia conosciuta e occorra stimarla, la formula dell'intervallo si modifica come segue

$$\hat{\mu} \pm t_{1-\alpha/2, n-1} \frac{\hat{\sigma}}{\sqrt{n}} \quad (26.4)$$

da cui

$$E = t_{1-\alpha/2, n-1} \frac{\sigma}{\sqrt{n}} \quad (26.5)$$

e specularmente per il calcolo dell'ampiezza:

$$n = \frac{t_{1-\alpha/2, n-1}^2 \cdot \hat{\sigma}^2}{E^2} \quad (26.6)$$

### 26.1.1.2 Intervallo a una coda

Nel caso di intervallo ad una coda, si prende come misura di errore massimo la distanza tra la stima e il valore dell'intervallo calcolato (non quello infinito, chiaramente). Ad esempio il limite inferiore di un intervallo di confidenza ad una coda (facciamo l'esempio di varianza nota per brevità di notazione)

$$L = \hat{\mu} - z_{1-\alpha} \frac{\sigma}{\sqrt{n}} \quad (26.7)$$

con l'intervallo che va da  $L$  a  $+\infty$ . L'errore massimo che vorremo commettere sarà dunque paria a

$$E = z_{1-\alpha} \frac{\sigma}{\sqrt{n}} \quad (26.8)$$

e risolvendo per  $n$  la formula per il dimensionamento è

$$n = \frac{z_{1-\alpha}^2 \cdot \sigma^2}{E^2} \quad (26.9)$$

**Example 26.1.2.** Per un intervallo ad una coda al 95% con errore non superiore al 10% della deviazione standard: si ha che che l'errore massimo è

$$z_{1-\alpha} \frac{\sigma}{\sqrt{n}} = 0.1\sigma$$

E pertanto

$$n = \frac{z_{1-\alpha} \sigma^2}{E^2} = \frac{(1.65)^2 \sigma^2}{(0.1\sigma)^2} = 272.2 \approx 273$$

## 26.1.2 Stima di una proporzione

### 26.1.2.1 Intervallo a due code

La formula dell'intervallo asintotico di una proporzione è

$$\hat{p} \pm z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}}$$

da cui l'errore

$$E = z_{1-\alpha/2} \sqrt{\frac{p(1-p)}{n}} \quad (26.10)$$

e la formula del campione

$$n = \frac{z_{1-\alpha/2}^2 \cdot p \cdot (1-p)}{E^2} \quad (26.11)$$

**Example 26.1.3.** Se si desidera nell'ipotesi che  $p = 0.5$  che l'intervallo di confidenza abbia una semiampiezza di 0.05 si ha

$$\frac{1.96^2 \cdot 0.5 \cdot (1-0.5)}{0.05^2} = 384.1 \approx 385$$

*Remark 92.* A parità di livello di significatività ed ampiezza massima di errore della stima, la numerosità dipende dalla variabilità  $p(1-p)$  al numeratore; cautelativamente, qualora non si abbia una minima idea di quale possa essere la prevalenza, adottare  $p = 0.5$

**Example 26.1.4.** In figura 26.1 la semi ampiezza di un intervallo di confidenza al 95% per un campione di 150 pazienti al variare della prevalenza

```
e <- function(p, alpha = 0.05, tails = 2, n = 150){
  z <- qnorm(1 - alpha/tails)
  z * sqrt( p * (1-p) / n )
}

ps <- seq(0.1, 0.9, by = 0.01)
plot(x = ps, y = e(p = ps), xlab = 'p', ylab = 'E', pch = NA, ylim = c(0, 0.1))
lines(x = ps, y = e(p = ps))
```

## 26.2 Power analysis

Dato che un errore del primo tipo è solitamente considerato più importante/grave, un approccio tipico al test di ipotesi è di controllare  $\alpha$  ad un livello accettabile e cercare di minimizzare  $\beta$  scegliendo un sample size adeguato.

### 26.2.1 Test per una media

Qui supponiamo di avere dati di una variabile quantitativa su  $n$  soggetti  $x_1, \dots, x_n$ ; media e varianza campionaria sono rispettivamente

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (26.12)$$

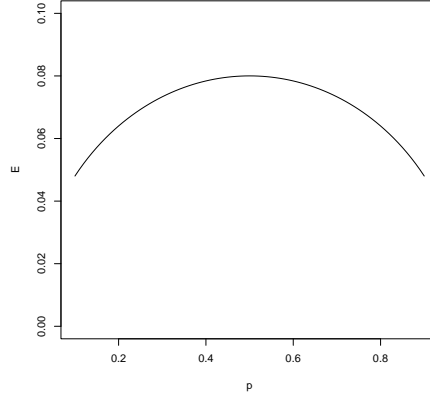


Figura 26.1: Semiampiezza intervallo di confidenza per 150 pz e varie prevalenze

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (26.13)$$

### 26.2.1.1 Equivalenza

Vogliamo verificare se la risposta media della popolazione  $\mu$  sia differente o meno da un valore di riferimento  $\mu_0$  e denominiamo  $\epsilon = \mu - \mu_0$  la differenza. Chiaramente:

$$\epsilon = 0 \iff \mu = \mu_0$$

Per verificare se vi sia una differenza tra la risposta media e il valore di riferimento le ipotesi poste a confronto sono:

$$H_0 : \epsilon = 0, \quad H_1 : \epsilon \neq 0$$

**Varianza conosciuta** Ipotizzando di conoscere la deviazione standard  $\sigma$  del carattere nella popolazione, il test da applicare è lo z:

$$z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$$

si ha che:

- sotto ipotesi nulla che  $\epsilon = 0$ ,  $z$  si distribuisce come  $N(\epsilon, 1) = N(0, 1)$ , da cui deriva che con test bilaterale rifiutiamo  $H_0$  se  $|z| > z_{1-\alpha/2}$ , e se  $\alpha = 0.05$  per  $|z| > 1.96$ ;
- sotto ipotesi alternativa  $z$  si distribuisce come  $N(\epsilon^*, 1)$  con

$$\epsilon^* = \frac{\epsilon}{\sigma/\sqrt{n}} = \frac{\epsilon\sqrt{n}}{\sigma}$$

con  $\epsilon \neq 0$ . La potenza di tale test è la probabilità di ottenere un risultato oltre la soglia rifiuto nell'ipotesi che sia vera l'alternativa, ossia:

$$\begin{aligned}\mathbb{P}(|N(\epsilon^*, 1)| > z_{1-\alpha/2}) &= \mathbb{P}(N(\epsilon^*, 1) > z_{1-\alpha/2}) + \mathbb{P}(N(\epsilon^*, 1) < -z_{1-\alpha/2}) \\ &= \mathbb{P}(N(0, 1) > z_{1-\alpha/2} - \epsilon^*) + \mathbb{P}(N(0, 1) < -z_{1-\alpha/2} - \epsilon^*) \\ &\stackrel{(1)}{=} \mathbb{P}(N(0, 1) < \epsilon^* - z_{1-\alpha/2}) + \mathbb{P}(N(0, 1) < -z_{1-\alpha/2} - \epsilon^*) \\ &= \Phi(\epsilon^* - z_{1-\alpha/2}) + \Phi(-\epsilon^* - z_{1-\alpha/2}) \\ &= \Phi\left(\frac{\epsilon\sqrt{n}}{\sigma} - z_{1-\alpha/2}\right) + \Phi\left(-\frac{\epsilon\sqrt{n}}{\sigma} - z_{1-\alpha/2}\right)\end{aligned}$$

dove in (1) abbiamo sfruttato la simmetria della normale. Ignorando una piccola parte di potenza ( $\leq \alpha/2$ ), possiamo dire che la potenza è approssimativamente pari a

$$\Phi\left(\frac{\epsilon\sqrt{n}}{\sigma} - z_{1-\alpha/2}\right) + \Phi\left(-\frac{\epsilon\sqrt{n}}{\sigma} - z_{1-\alpha/2}\right) \approx \Phi\left(\left|\frac{\epsilon\sqrt{n}}{\sigma}\right| - z_{1-\alpha/2}\right) = \Phi\left(\frac{|\epsilon|\sqrt{n}}{\sigma} - z_{1-\alpha/2}\right)$$

Affinché la potenza sia  $1 - \beta$  con  $\beta$  scelto a piacere, si deve avere

$$\Phi\left(\frac{|\epsilon|\sqrt{n}}{\sigma} - z_{1-\alpha/2}\right) = 1 - \beta$$

ossia

$$\frac{|\epsilon|\sqrt{n}}{\sigma} - z_{1-\alpha/2} = z_{1-\beta}$$

e risolvendo per  $n$  si giunge a

$$\begin{aligned}\sqrt{n}|\epsilon| &= (z_{1-\beta} + z_{1-\alpha/2})\sigma \\ n &= \frac{(z_{1-\beta} + z_{1-\alpha/2})^2\sigma^2}{\epsilon^2}\end{aligned}$$

In merito a quest'ultima alcuni libri scrivono  $z_\beta$  al posto di  $z_{1-\beta}$  e  $z_{\alpha/2}$  al posto di  $z_{1-\alpha/2}$ , ma è corretto come si è fatto qui.

**Example 26.2.1.** Riproducendo, se  $\alpha = 0.05$ ,  $\beta = 0.2$ ,  $\sigma = 1$ ,  $\epsilon = 0.5$  si ha

```
num <- (qnorm(1-0.05/2) + qnorm(1 - 0.2))^2 * 1^2
den <- 0.5^2
num / den

## [1] 31.39552
```

**Varianza ignota** Quando  $\sigma^2$  è sconosciuta può esser rimpiazzata dalla varianza campionaria data in 26.13 e l'ipotesi  $H_0$  viene rifiutata se

$$\left|\frac{\bar{x} - \mu_0}{s/\sqrt{n}}\right| > t_{1-\alpha/2, n-1}$$

dove  $t_{1-\alpha/2, n-1}$  è il quantile  $1 - \alpha/2$  della distribuzione  $t$  con  $n - 1$  gradi di libertà

```

beta <- function(theta, alpha = 0.05, max_beta = 0.2, N = 2:45){
  # per ogni campione indagato, fai il calcolo di beta ottenuto
  res <- lapply(N, function(n){
    df <- n - 1
    ncp <- sqrt(n) * theta
    t <- qt(1 - alpha/2, df = df, ncp = 0)
    beta <- pt(q = t, df = n-1, ncp = ncp)
    data.frame(n = n, beta = beta, ok = beta < max_beta)
  })
  do.call(rbind, res)
}

beta(theta = 0.5)

```

##	n	beta	ok
## 1	2	0.94690414	FALSE
## 2	3	0.92106852	FALSE
## 3	4	0.89191676	FALSE
## 4	5	0.86154723	FALSE
## 5	6	0.83067614	FALSE
## 6	7	0.79962458	FALSE
## 7	8	0.76859262	FALSE
## 8	9	0.73773267	FALSE
## 9	10	0.70717143	FALSE
## 10	11	0.67701744	FALSE
## 11	12	0.64736426	FALSE
## 12	13	0.61829219	FALSE
## 13	14	0.58986950	FALSE
## 14	15	0.56215341	FALSE
## 15	16	0.53519106	FALSE
## 16	17	0.50902038	FALSE
## 17	18	0.48367096	FALSE
## 18	19	0.45916479	FALSE
## 19	20	0.43551710	FALSE
## 20	21	0.41273703	FALSE
## 21	22	0.39082828	FALSE
## 22	23	0.36978977	FALSE
## 23	24	0.34961622	FALSE
## 24	25	0.33029864	FALSE
## 25	26	0.31182487	FALSE
## 26	27	0.29418002	FALSE
## 27	28	0.27734685	FALSE
## 28	29	0.26130620	FALSE
## 29	30	0.24603728	FALSE
## 30	31	0.23151800	FALSE
## 31	32	0.21772526	FALSE
## 32	33	0.20463518	FALSE
## 33	34	0.19222331	TRUE
## 34	35	0.18046489	TRUE
## 35	36	0.16933492	TRUE

```
## 36 37 0.15880842 TRUE
## 37 38 0.14886050 TRUE
## 38 39 0.13946650 TRUE
## 39 40 0.13060209 TRUE
## 40 41 0.12224334 TRUE
## 41 42 0.11436681 TRUE
## 42 43 0.10694963 TRUE
## 43 44 0.09996950 TRUE
## 44 45 0.09340477 TRUE
```

Per cui si nota che all'aumentare del campione l'errore di secondo tipo diminuisce e ad un sample size di 34 si ottiene un  $\beta < 0.2$  mentre per una potenza del 90% occorrono 44 soggetti.

### 26.2.1.2 Superiority/Non-inferiority

Disegni di non inferiorità e superiorità possono essere unificati dal seguente test a una coda

$$H_0 : \epsilon \leq \delta, \quad H_1 : \epsilon > \delta$$

con  $\delta$  detto margine di superiorità o non inferiorità. Quando:

- $\delta = 0$ : il disegno è una superiorità classica (in senso statistico), ossia confronta  $H_0 : \mu \leq \mu_0$  con  $H_1 : \mu > \mu_0$
- $\delta > 0$ : disegno di superiorità che verifica che la differenza tra media della popolazione e valore teorico sia superiore ad un dato valore  $\delta$
- $\delta < 0$ : il disegno è di non inferiorità; mira a verificare che il valore della popolazione sia entro una certa distanza  $\delta$  dal valore ipotizzato  $\mu_0$ .  
In altre parole se si verifica la nulla  $\mu - \mu_0 \leq \delta$ , ossia  $\mu \leq \mu_0 + \delta$  (con  $\delta < 0$ ); se si verifica l'alternativa  $\mu - \mu_0 > \delta$  ossia  $\mu_0 - \mu < -\delta$  con  $-\delta > 0$

**Varianza nota** Ipotizzando che  $\sigma^2$  sia nota il test che impieghiamo per la scelta è

$$\frac{\bar{x} - \mu_0 - \delta}{\sigma/\sqrt{n}}$$

Pertanto è come se si passasse dal test sul valore puntuale vs valore teorico al test della differenza  $(\bar{x} - \mu_0)$  vs una sorta di valore teorico della differenza  $\delta$ . Il test rispetto ad una  $z$  normale non cambia (perché è come se incrementassimo la costante data dall'ipotesi nulla al numeratore), quindi si distribuisce sempre con  $N(0, 1)$ . Essendo un test ad una coda sul valore scarto  $\epsilon$  vs  $\delta$  la zona di rifiuto di  $H_0$  è tutta a destra, per cui  $H_0$  è rifiutata se

$$\frac{\bar{x} - \mu_0 - \delta}{\sigma/\sqrt{n}} > z_{1-\alpha}$$

Il test ha una distribuzione normale centrata sul valore  $\delta$ ; sotto ipotesi alternativa, se  $\epsilon > \delta$  la distribuzione dello stimatore non è centrata su  $\delta$  ma su  $\epsilon - \delta$ , che in unità di misura della distribuzione nulla sono

$$\frac{\epsilon - \delta}{\sigma/\sqrt{n}}$$

quindi sotto ipotesi alternativa il test si distribuisce come

$$z \sim N\left(\frac{\epsilon - \delta}{\sigma/\sqrt{n}}, 1\right)$$

Siamo interessati alla potenza del test, ossia a:

$$\begin{aligned} \mathbb{P}\left(N\left(\frac{\epsilon - \delta}{\sigma/\sqrt{n}}, 1\right) > z_{1-\alpha}\right) &= \mathbb{P}\left(N(0, 1) > z_{1-\alpha} - \frac{\epsilon - \delta}{\sigma/\sqrt{n}}\right) \\ &= \mathbb{P}\left(N(0, 1) < \frac{\epsilon - \delta}{\sigma/\sqrt{n}} - z_{1-\alpha}\right) \\ &= \Phi\left(\frac{\epsilon - \delta}{\sigma/\sqrt{n}} - z_{1-\alpha}\right) \end{aligned}$$

Impostandola al livello di potenza desiderato:

$$\Phi\left(\frac{\epsilon - \delta}{\sigma/\sqrt{n}} - z_{1-\alpha}\right) = 1 - \beta$$

da cui

$$\frac{\epsilon - \delta}{\sigma/\sqrt{n}} - z_{1-\alpha} = z_{1-\beta}$$

e risolvendo per  $n$  si giunge a

$$n = \frac{(z_{1-\alpha} + z_{1-\beta})^2 \sigma^2}{(\epsilon - \delta)^2}$$

**Varianza sconosciuta**

### 26.2.1.3 Equivalence

### 26.2.2 Test per una proporzione



## Capitolo 27

# Two groups

### 27.1 T-test

Supposing<sup>1</sup> the endpoint for a comparative trial is a quantitative measurement so the treatment comparison consists in testing the difference of estimated mean of the two treatment groups.

Supposing :

- the true means in the treatment groups are  $\mu_t$  and  $\mu_c$  and the common standard deviation for the two groups in the population is  $\sigma$ .
- the treatment difference will be  $\Delta = \mu_t - \mu_c$
- the hypothesis compared will be

$$\begin{aligned} H_0 : \Delta = 0 & \iff \mu_t = \mu_c \\ H_1 : \Delta \neq 0 & \iff \mu_t \neq \mu_c \end{aligned}$$

- the investigator would reject the null hypothesis if the observed  $|\Delta|$  exceeds a critical value  $c$  that is (under normal distribution of the estimator mean difference)

$$c = z_{1-\alpha/2} \cdot \sigma_\Delta$$

where  $\sigma_\Delta$  is the sd of the estimator mean difference and  $z_{1-\alpha/2}$  is the  $1 - \alpha/2$  quantile of the standard normal distribution

Now we want to express such  $c$  as quantile of the distribution under alternative hypotheses. The estimator has a normal distribution which under the alternative is centered on  $\Delta$  so to express the rejection cutoff as quantile/z score in this distribution we have

$$z_\beta = \frac{z_{1-\alpha/2} \cdot \sigma_\Delta - \Delta}{\sigma_\Delta} = z_{1-\alpha/2} - \frac{\Delta}{\sigma_\Delta}$$

So it turns out that the standardized effect size is defined as

$$\frac{\Delta}{\sigma_\Delta} = z_{1-\alpha/2} - z_\beta = z_{1-\alpha/2} + z_{1-\beta} \quad (27.1)$$

---

<sup>1</sup>Da [piantadosi2005rct], modificata notazione

where the symmetry property of normal was used. Now from statistical theory it turns out that the difference between two means (the estimator) has standard deviation

$$\sigma_{\Delta} = \sqrt{\frac{\sigma^2}{n_c} + \frac{\sigma^2}{n_t}} = \sigma \sqrt{\frac{1}{n_c} + \frac{1}{n_t}}$$

where we assumed a common standard deviation for the two groups in the population  $\sigma_c = \sigma_t = \sigma$ . Substituting this quantity in 27.1 and squaring both terms we obtain

$$\begin{aligned} z_{1-\alpha/2} + z_{1-\beta} &= \frac{\Delta}{\sigma \sqrt{\frac{1}{n_c} + \frac{1}{n_t}}} \\ (z_{1-\alpha/2} + z_{1-\beta})^2 &= \frac{\Delta^2}{\sigma^2 \cdot \left(\frac{1}{n_c} + \frac{1}{n_t}\right)} \\ \left(\frac{1}{n_c} + \frac{1}{n_t}\right) &= \frac{\Delta^2}{\sigma^2 (z_{1-\alpha/2} + z_{1-\beta})^2} \end{aligned}$$

Now if we set the allocation ratio  $r = \frac{n_t}{n_c}$  so that  $n_t = r \cdot n_c$  we have

$$\frac{1}{n_c} + \frac{1}{r \cdot n_c} = \frac{1}{n_c} \cdot \frac{r+1}{r} = \frac{\Delta^2}{\sigma^2 (z_{1-\alpha/2} + z_{1-\beta})^2}$$

Thus once fixed  $r$ ,  $n_c$  can be obtained as

$$n_c = \frac{r+1}{r} \cdot \frac{\sigma^2 (z_{1-\alpha/2} + z_{1-\beta})^2}{\Delta^2}$$

If  $r = 1$  and thus  $n_c = n_t$ , then  $n = n_c + n_t$  will be

$$\begin{aligned} n = n_c + n_t &= 2 \cdot n_c = 2 \cdot \frac{1+1}{1} \cdot \frac{\sigma^2 (z_{1-\alpha/2} + z_{1-\beta})^2}{\Delta^2} \\ &= \frac{4\sigma^2 (z_{1-\alpha/2} + z_{1-\beta})^2}{\Delta^2} \end{aligned}$$

This formula allows us to calculate the sample size once some properties of the test (alpha, beta, effect size, common standard deviation of the quantitative variable in the two group) is given.

Otherwise we could determine the power in function of other parameters which is the integral

$$1 - \Phi(z_{\beta}) = 1 - \beta$$

To find  $z_\beta$  from above we start from the sample size formula and rearrange a bit and take the root to both terms

$$n = \frac{4\sigma^2}{\Delta^2}$$

$$(z_{1-\alpha/2} + z_{1-\beta})^2 = \frac{n\Delta^2}{4\sigma^2}$$

$$(z_{1-\alpha/2} + z_{1-\beta}) = \frac{\sqrt{n}|\Delta|}{2\sigma}$$

$$z_{1-\beta} = \frac{\sqrt{n} \cdot |\Delta|}{2\sigma} - z_{1-\alpha/2}$$

thus being interested in  $1 - \beta$  the power is (taking  $\Phi$  on both terms)

$$1 - \beta = \Phi\left(\frac{\sqrt{n} \cdot |\Delta|}{2\sigma} - z_{1-\alpha/2}\right)$$

Some remarks:

- this power equation has been derived assuming normal distribution of the test statistics and assuming a known common standard deviation  $\sigma$  in the population
- the CLT assures us that the distribution of the test statistics will be approximately normal (if groups size are not too scarce which is typically not the case in trials) while for known variance, if variance have to be estimated by the data the distribution of the test statistics will be approximately be a T distribution (not normal anymore) but the difference in sample size produced is negligible as sample size increases

```
# check differenze per vari effect size standardizzati sotto queste ipotesi
std_es <- seq(0.05, 1.5, by = 0.05)
alpha <- 0.05
beta <- 0.2

res <- lapply(std_es, function(es){
  t_nc <- power.t.test(delta = es, power = 1-beta, sig.level = alpha)$n
  z_nc <- 2 * 1 * (qnorm(1 - alpha/2) + qnorm(1 - beta))^2 / (es^2)
  data.frame("Effect size"= es,
             "ss_z" = ceiling(z_nc) * 2,
             "ss_t" = ceiling(t_nc) * 2)
})

do.call(rbind, res)

##      Effect.size  ss_z  ss_t
## 1          0.05 12560 12562
## 2          0.10  3140  3142
## 3          0.15  1396  1398
## 4          0.20   786   788
## 5          0.25   504   506
```

## 6	0.30	350	352
## 7	0.35	258	260
## 8	0.40	198	200
## 9	0.45	156	158
## 10	0.50	126	128
## 11	0.55	104	106
## 12	0.60	88	90
## 13	0.65	76	78
## 14	0.70	66	68
## 15	0.75	56	58
## 16	0.80	50	52
## 17	0.85	44	46
## 18	0.90	40	42
## 19	0.95	36	38
## 20	1.00	32	34
## 21	1.05	30	32
## 22	1.10	26	30
## 23	1.15	24	26
## 24	1.20	22	24
## 25	1.25	22	24
## 26	1.30	20	22
## 27	1.35	18	20
## 28	1.40	18	20
## 29	1.45	16	18
## 30	1.50	14	18

*# la differenza diminuisce mano a mano. a parte i numeri bassi rimane ignorabile*

## Capitolo 28

# Multiple endpoints

### 28.1 Methodology

**Application** Following FDA ([16]), multiple endpoints may be **needed when**, to determine that a treatment confers benefit, when either:

- there are several important aspects of a disease
- there are several ways to assess a same important aspect

and either

- there is no consensus about which one will best serve the study purposes
- there's no way to combine the 2+ outcomes, or it would be considered an information lost

**Types of multiple outcome studies** If multiple endpoints are used then, different strategies can be adopted to **state that a treatment** is effective:

- for some disease, there are two or more features that are so critically important that a treatment will not be considered effective without demonstration of a treatment effect on *all* of these disease features. In this case multiple outcomes are called *co-primary* outcomes, and it is necessary to demonstrate an effect on each of the endpoints to conclude that a treatment is effective.
- in other cases a study might be designed such that success on any one of several endpoints is sufficient to support a conclusion of effectiveness (*multiple primary* outcomes).

**Type 1 Errors** If the study is with

- *co-primary outcomes*: in this case ([16] pag 13) there's no need to adjust first type error inflation in applying the p-value cutoff in each analysis.

Prove: if we reject the null if  $p < \alpha = 0.05$  and considering two independent outcome the overall type 1 error, that is probability of having two statistical significant results under null hypotheses is:

$$\alpha^* = \alpha^2 = 0.05^2$$

So there's no issue with type 1 errors and no need to change nominal level in each analysis to have overall type error below a certain threshold.

In other words there is non need for multiplicity adjustment because there is no opportunity to select the most favorable result from among several endpoint.

However the impact of multiplicity in these situations is to increase the type II error rate

- *multiple primary* outcomes: success on any one of several endpoints is sufficient to support a conclusion of effectiveness.

In this case there's an inflation of Type I error (over the canonical two-sided 0.05 or one sided 0.025) if both:

- each separate outcome analysis is conducted at the default levels: eg to declare *statistical significant result, on each outcome*, we use two-sided 0.05 or one sided 0.025;
- to declare the *treatment effective* if any of the outcome analysis has statistical significant results (as is with multiple primary).

```
inflated_05 <- 1 - (1 - 0.05)^2
inflated_025 <- 1 - (1 - 0.025)^2
```

Eg: in case of two independent outcomes/analyses the overall Type I error (that is the probability of declaring the treatment effective, while it is not on both outcomes, if any of the outcome analyses is statistically significant) inflates from 5% to 9.75% (bilateral) and from 0.025 to 4.9375%.

Various procedures can be used to keep the overall Type I error at the desired phase.

**Type II and power** As the FDA Guidance for Industry points out [16], multiple outcomes changes the global Type II error (and power).

This is specular to what occurs for type I error, that is for study with

- *co-primary outcomes*: overall type II error  $\beta$  increases (and power decrease). In terms of power, in the case of two endpoints

$$\begin{aligned} 1 - \beta^{FWER} &= \mathbb{P}(\text{Refuse1}|H_{11}) \cap \mathbb{P}(\text{Refuse2}|H_{12}) \\ &= (1 - \beta)^2 \\ \beta^{FWER} &= 1 - (1 - \beta)^2 \end{aligned}$$

eg if the study sample size is able to provide 80% power ( $\beta = 0.2$ ) to show success on each endpoint separately, and endpoints are independent, the probability of rejecting both the null hypotheses under alternative hypotheses for both cases is  $0.8^2 = 0.64$  and thus the type II error (probability

of failing in rejecting at least a wrong null hypotheses) becomes 36%.  
 The loss of power may not be so severe when endpoints are correlated/no more independent, but have to be considered nonetheless. Thus sample size determination has to be considered for multiplicity of outcomes.

- *multiple-primary outcomes*: if it is enough to have a single statistical significant result to determine treatment efficacy the power is the probability of having

$$\begin{aligned} 1 - \beta^{FWER} &= \mathbb{P}(Refuse1|H_{11}) \cup \mathbb{P}(Refuse2|H_{12}) \\ &= 1 - \mathbb{P}(NotRefuse1|H_{11}) \cup \mathbb{P}(NotRefuse2|H_{12})\beta^{FWER} = \beta^2 \end{aligned}$$

so here actually beta decrease and power increase. However in this case of studies we have to control type I error which affect beta and power as well so sample size consideration are needed as well

**Summing up** Following [23] utilizing multiple endpoints may provide the opportunity for characterizing interventions's multidimensional effects, but also creates challenges.

Specifically controlling type I and type II error rates is non-trivial when the multiple primary endpoints are potentially correlated.

When more than one endpoint is viewed as important in a clinical trial, then a decision must be made as to whether it is desirable to evaluate the joint effects on ALL endpoints or AT LEAST ONE of the endpoints.

This decision *defines the alternative hypothesis to be tested* and provides a framework for trial design:

- when designing the trial to evaluate the joint effects on ALL of the endpoints, no adjustment is needed to control the type I error rate. The hypothesis associated with each endpoint can be evaluated at the same significance level that is desired for demonstrating effects on all endpoints.  
 However the type 2 error rate increases as the number of endpoints to be evaluated and power diminish
- when designing trial to evaluate an effect on AT LEAST ONE of the endpoints, an adjustment is needed to control the type I error rate. In this trials the **correlation** among multiple endpoints should be considered in order to obtain an appropriate sample size

## 28.2 MPE

Il pacchetto `mpe` quando c'era forniva due funzioni di interesse: `power.known.var` (multiple co-primary) e `atleast.one.endpoint` (multiple primary).

Da attenzione che `atleast.one.endpoint` pone di default il test a una coda a 5/numero confronti non a 2.5/numero confronti come dovrebbe essere per avere in ciascun test l'equivalentedi un test a due code al 5,

Si è estratta la parte importante e posta in `lbss`

```
library(lbss)

## Error in library(lbss): non c'è alcun pacchetto chiamato 'lbss'

## Esempio 2 continuous coprimary endpoints
example(me_both_the_two)

## Warning in example(me_both_the_two): nessun aiuto per 'me_both_the_two'

## Esempio 2 continuous primary endpoints
example(me_atleastone_among_two)

## Warning in example(me_atleastone_among_two): nessun aiuto per 'me_atleastone_amon
```



# Bibliografia

- [1] Douglas G. Altman. “Comparability of Randomised Groups”. In: *The Statistician* 34.1 (1985). Publisher: JSTOR, p. 125. DOI: 10.2307/2987510. URL: <https://doi.org/10.2307/2987510>.
- [2] Dino Amadori. *Sperimentazione clinica in oncologia*. Poletto Editore, 2004.
- [3] P. Armitage. “Tests for linear trends in proportions and frequencies”. In: *Biometrics* 11.3 (set. 1955), pp. 375–386.
- [4] Antonella Bacchieri e Giovanni Della Cioppa. *Fondamenti di ricerca clinica*. ita. Medicina e statistica. Milano Berlin Heidelberg: Springer, 2004. ISBN: 978-88-470-0211-1.
- [5] Gian Luca Baio. *Bayesian methods in health economics*. eng. Chapman & Hall/CRC biostatistics series 53. Boca Raton, Fla: Chapman & Hall/CRC, 2013. ISBN: 978-1-4398-9556-6.
- [6] Michael Borenstein, cur. *Introduction to meta-analysis*. eng. Statistics in Practice. Chichester, U.K: John Wiley & Sons, 2009. ISBN: 978-1-119-96437-7 978-0-470-74338-6 978-0-470-74337-9.
- [7] Sarah R. Brown et al., cur. *A practical guide to designing phase II trials in oncology*. Chichester, West Sussex: John Wiley & Sons Inc, 2014. ISBN: 978-1-118-76363-6 978-1-118-76362-9.
- [8] John Bryant e Roger Day. “Incorporating Toxicity Considerations Into the Design of Two-Stage Phase II Clinical Trials”. In: *Biometrics* 51.4 (dic. 1995), p. 1372. ISSN: 0006341X. DOI: 10.2307/2533268. URL: <https://www.jstor.org/stable/2533268?origin=crossref> (visitato il giorno 01/08/2025).
- [9] Myron N. Chang et al. “Designs for Group Sequential Phase II Clinical Trials”. In: *Biometrics* 43.4 (dic. 1987), p. 865. ISSN: 0006341X. DOI: 10.2307/2531540. URL: <https://www.jstor.org/stable/2531540?origin=crossref> (visitato il giorno 31/07/2025).
- [10] T. Timothy Chen. “Optimal three-stage designs for phase II cancer clinical trials”. en. In: *Statistics in Medicine* 16.23 (dic. 1997), pp. 2701–2711. ISSN: 0277-6715, 1097-0258. DOI: 10.1002/(SICI)1097-0258(19971215)16:23<2701::AID-SIM704>3.0.CO;2-1. URL: [https://onlinelibrary.wiley.com/doi/10.1002/\(SICI\)1097-0258\(19971215\)16:23%3C2701::AID-SIM704%3E3.0.CO;2-1](https://onlinelibrary.wiley.com/doi/10.1002/(SICI)1097-0258(19971215)16:23%3C2701::AID-SIM704%3E3.0.CO;2-1) (visitato il giorno 01/08/2025).

- [11] Shein-Chung Chow, Jun Shao e Hansheng Wang. *Sample Size Calculations in Clinical Research*. eng. 2nd ed. Chapman & Hall/CRC biostatistics series 20. OCLC: 890665928. Boca Raton: Chapman & Hall/CRC, 2008. ISBN: 978-1-58488-982-3.
- [12] E.A. Eisenhauer et al. "New response evaluation criteria in solid tumours: Revised RECIST guideline (version 1.1)". en. In: *European Journal of Cancer* 45.2 (gen. 2009), pp. 228–247. ISSN: 09598049. DOI: 10.1016/j.ejca.2008.10.026. URL: <https://linkinghub.elsevier.com/retrieve/pii/S0959804908008733> (visitato il giorno 01/08/2025).
- [13] Elizabeth A. Eisenhauer. *Phase I Cancer Clinical Trials: A Practical Guide (2nd Edition)*. eng. 2nd ed. New York: Oxford University Press, Incorporated, 2015. ISBN: 978-0-19-935901-1 978-0-19-935902-8.
- [14] S. M. Eldridge et al. "CONSORT 2010 statement: extension to randomised pilot and feasibility trials". In: *BMJ* 355 (ott. 2016), p. i5239.
- [15] Lisa Garnsey Ensign et al. "An optimal three-stage design for phase II clinical trials". en. In: *Statistics in Medicine* 13.17 (set. 1994), pp. 1727–1736. ISSN: 0277-6715, 1097-0258. DOI: 10.1002/sim.4780131704. URL: <https://onlinelibrary.wiley.com/doi/10.1002/sim.4780131704> (visitato il giorno 31/07/2025).
- [16] FDA. *Multiple Endpoints in Clinical Trials*. en. Publisher: FDA. Ago. 2022. URL: <https://www.fda.gov/regulatory-information/search-fda-guidance-documents/multiple-endpoints-clinical-trials> (visitato il giorno 05/08/2025).
- [17] Henry Glick, cur. *Economic evaluation in clinical trials*. Handbooks in health economic evaluation series. Oxford ; New York: Oxford University Press, 2007. ISBN: 978-0-19-852997-2.
- [18] Julian P.T. Higgins et al., cur. *Cochrane Handbook for Systematic Reviews of Interventions*. en. 1<sup>a</sup> ed. Wiley, set. 2019. ISBN: 978-1-119-53662-8 978-1-119-53660-4. DOI: 10.1002/9781119536604. URL: <https://onlinelibrary.wiley.com/doi/book/10.1002/9781119536604> (visitato il giorno 31/07/2025).
- [19] ICMJE. *Recommendations for the Conduct, Reporting, Editing, and Publication of Scholarly Work in Medical Journals*. Dic. 2014.
- [20] L. Mariani e E. Marubini. "Design and Analysis of Phase II Cancer Trials: A Review of Statistical Methods and Guidelines for Medical Researchers". In: *International Statistical Review / Revue Internationale de Statistique* 64.1 (apr. 1996), p. 61. ISSN: 03067734. DOI: 10.2307/1403424. URL: <https://www.jstor.org/stable/1403424?origin=crossref> (visitato il giorno 31/07/2025).
- [21] Margaret Sullivan Pepe. *The Statistical Evaluation of Medical Tests for Classification and Prediction*. en. Oxford University Press Oxford, mar. 2003. ISBN: 978-0-19-850984-4 978-1-383-02220-9. DOI: 10.1093/oso/9780198509844.001.0001. URL: <https://academic.oup.com/book/52788> (visitato il giorno 01/08/2025).

- [22] Richard Simon. “Optimal two-stage designs for phase II clinical trials”. en. In: *Controlled Clinical Trials* 10.1 (mar. 1989), pp. 1–10. ISSN: 01972456. DOI: 10.1016/0197-2456(89)90015-9. URL: <https://linkinghub.elsevier.com/retrieve/pii/0197245689900159> (visitato il giorno 01/08/2025).
- [23] Takashi Sozu et al. *Sample Size Determination in Clinical Trials with Multiple Endpoints*. en. SpringerBriefs in Statistics. Cham: Springer International Publishing, 2015. ISBN: 978-3-319-22004-8 978-3-319-22005-5. DOI: 10.1007/978-3-319-22005-5. URL: <https://link.springer.com/10.1007/978-3-319-22005-5> (visitato il giorno 01/08/2025).
- [24] Graham M. Wheeler et al. “How to design a dose-finding study using the continual reassessment method”. In: *BMC Medical Research Methodology* 19.1 (gen. 2019). Publisher: Springer Science and Business Media LLC. DOI: 10.1186/s12874-018-0638-z. URL: <https://doi.org/10.1186/s12874-018-0638-z>.
- [25] Mark Woodward. *Epidemiology: study design and data analysis*. eng. 2. ed. Chapman & Hall/CRC texts in statistical science series. Boca Raton, Fla. London: Chapman & Hall/CRC, 2005. ISBN: 978-1-58488-415-6.