

Indice

Prefazione	11
1. Obiettivi, modelli e metodi della ricerca sociale	17
1.1. I metodi della ricerca	18
1.1.1. <i>Il principio di falsificazione</i>	18
1.1.2. <i>Evoluzione del metodo scientifico</i>	19
1.1.3. <i>Metodologia della ricerca</i>	20
1.2. Tipi di ricerca	21
1.3. Disegno della ricerca	23
2. Strumenti qualitativi	27
2.1. Un quadro d'insieme	27
2.2. Ricerca qualitativa e quantitativa	28
2.3. Metodi di ricerca qualitativa	31
2.3.1. <i>Studi osservazionali</i>	31
2.3.2. <i>L'intervista qualitativa</i>	34
2.3.3. <i>Tecniche di gruppo</i>	36
2.3.4. <i>Ricerca documentale</i>	40
2.4. Intervista qualitativa	42
2.5. Osservazione ed osservazione partecipata	46
2.6. Focus group	48
2.6.1. <i>Modalità operative</i>	50
2.6.2. <i>Documenti del focus group – Lettera di presentazione</i>	53
2.6.3. <i>Documenti del focus group – Domande per la discussione</i>	55
2.6.4. <i>Documenti del focus group – Questionario finale</i>	56
2.6.5. <i>Relazione sui dati raccolti</i>	59
2.7. Intervista a testimoni privilegiati	61
2.7.1. <i>Il metodo Delphi</i>	61
2.7.2. <i>Il metodo Delphi – Shang</i>	64

2.8.	Analisi di dati qualitativi	65
2.8.1.	<i>Carta tematica</i>	66
2.8.2.	<i>Analisi del contesto KWIC</i>	67
2.8.3.	<i>Vocabolario e distribuzione di frequenze</i>	69
2.8.4.	<i>Matrice parole per documento TDM</i>	71
2.8.5.	<i>Rango e frequenze</i>	73
2.8.6.	<i>Word Cloud</i>	74
2.8.7.	<i>Distanze testuali</i>	75
2.8.8.	<i>Vettori e Matrici</i>	79
2.8.9.	<i>Distanza del coseno con Excel</i>	83
2.8.10.	<i>Analisi dei gruppi</i>	84
2.8.11.	<i>Cluster analysis: metodo del legame medio</i>	87
2.8.12.	<i>Cluster Analysis con R</i>	89
3.	Strumenti quantitativi	101
3.1.	Fonti statistiche ufficiali	101
3.1.1.	<i>Stimatori per piccole aree</i>	103
3.1.2.	<i>Stima indiretta dell'occupazione a livello comunale</i>	104
3.2.	I dati documentali	110
3.2.1.	<i>I dati amministrativi</i>	111
3.3.	Inchiesta campionaria	112
3.3.1.	<i>Indagine faccia a faccia</i>	115
3.3.2.	<i>Indagine telefonica</i>	117
3.3.3.	<i>Indagine con questionario autocompilato</i>	120
3.4.	Costruzione del campione	124
3.4.1.	<i>Principio del campionamento ripetuto</i>	126
3.4.2.	<i>Errore di stima</i>	127
3.4.3.	<i>Teorema del limite centrale</i>	128
3.4.4.	<i>Varianza del campione e Varianza della popolazione</i>	130
3.4.5.	<i>Numerosità campionaria</i>	132
3.4.6.	<i>Esempio 1</i>	136
3.4.7.	<i>Esempio 2</i>	139
3.5.	Disegni di campionamento probabilistico	139
3.5.1.	<i>Campionamento casuale semplice</i>	140
3.5.2.	<i>Campionamento stratificato</i>	141
3.5.3.	<i>Campionamento a grappoli</i>	142
3.5.4.	<i>Campionamento a due o più stadi</i>	142
3.5.5.	<i>Campionamento sistematico</i>	143
3.5.6.	<i>Indagine longitudinale – Panel</i>	144
3.5.7.	<i>Indagine longitudinale – Panel ruotato</i>	144
3.5.8.	<i>Campionamento di aree</i>	145

3.6.	Disegni di campionamento non probabilistico	146
3.6.1.	<i>Campionamento ragionato</i>	146
3.6.2.	<i>Campionamento per quote</i>	146
3.6.3.	<i>Campionamento a valanga</i>	147
4.	Il questionario	149
4.1.	Caratteristiche generali	151
4.2.	Formulazione dei quesiti	152
4.2.1.	<i>Linguaggio utilizzato</i>	152
4.2.2.	<i>Categorie di domande</i>	153
4.3.	Grafica ed impaginazione	157
4.3.1.	<i>Lettera di presentazione</i>	157
4.3.2.	<i>Prima pagina del questionario</i>	158
4.3.3.	<i>Pagine successive</i>	159
4.3.4.	<i>Pagina conclusiva</i>	160
4.3.5.	<i>Codifica delle modalità di risposta</i>	160
4.4.	Le scale di misurazione	161
4.4.1.	<i>Variabili qualitative su scala nominale</i>	162
4.4.2.	<i>Variabili qualitative su scala ordinale</i>	163
4.4.3.	<i>Variabili quantitative su scala ad intervallo</i>	165
4.4.4.	<i>Variabili quantitative su scala a rapporto</i>	165
4.4.5.	<i>Le scale di misurazione secondo Ricolfi</i>	166
4.5.	Misura degli atteggiamenti	168
4.5.1.	<i>Tecniche di scaling</i>	169
4.6.	Scale comparative	170
4.6.1.	<i>Alcune scale di confronto</i>	171
4.6.2.	<i>Scale a somma costante</i>	177
4.6.3.	<i>Scale Q – sort o ordinamento qualitativo</i>	178
4.7.	Scale non comparative	180
4.7.1.	<i>Scala semanticamente autonoma</i>	182
4.7.2.	<i>Scale a parziale autonomia semantica</i>	182
4.7.3.	<i>Scale autoancoranti</i>	182
4.7.4.	<i>Scala di Likert</i>	183
4.7.5.	<i>Scala del differenziale semantico</i>	187
4.7.6.	<i>Scala di Stapel</i>	187
4.7.7.	<i>Scala di Guttman</i>	188
4.7.8.	<i>Scala di Thurstone</i>	189
4.8.	Batterie di domande	191

5. Analisi del Rischio	193
5.1. Studi trasversali e studi longitudinali	194
5.2. Alcune definizioni	194
5.3. Rischio per casi dicotomici	196
5.4. Alcuni esempi	199
5.4.1. <i>Inferenza sul rischio relativo</i>	200
5.4.2. <i>Inferenza sul rapporto tra Odds</i>	202
5.5. Esercizi	205
5.6. Analisi del rischio in R	212
6. Confronto fra campioni	221
6.1. Test t di Student per confronto tra due campioni indipendenti	221
6.1.1. <i>Varianze note</i>	222
6.1.2. <i>Varianze incognite ma uguali</i>	225
6.1.3. <i>Varianze incognite e diverse</i>	227
6.2. Test su proporzioni	230
6.2.1. <i>Test asintotico Z per una proporzione</i>	231
6.2.2. <i>Test su due proporzioni con varianza non nota ma uguale</i>	232
6.2.3. <i>Intervallo di confidenza per una proporzione</i>	233
7. Analisi di regressione multivariata	235
7.1. Regressione lineare multivariata	235
7.2. Regressione logistica	240
7.2.1. <i>Esempio: programma di riabilitazione cardiaca</i>	242
7.3. Interpretazione dei risultati	244
7.3.1. <i>Esercizio: infertilità dopo aborto spontaneo e aborto indotto</i>	245
7.3.2. <i>Esercizio: fecondità ed indicatori socio-economici</i>	248
7.3.3. <i>Esercizio: customer satisfaction</i>	249
8. Introduzione al linguaggio R	251
8.1. Introduzione	251
8.2. Assegnazioni e variabili in R	253
8.3. Inserimento di dati in una variabile	255
8.3.1. <i>Vettori</i>	257
8.3.2. <i>Matrici</i>	260
8.3.3. <i>Liste</i>	260
8.4. Trasformazione in ranghi per una variabile	261
8.5. Rappresentazioni grafiche: istogramma	262
8.5.1. <i>Varianti per il Grafico ad Istogramma</i>	263
8.6. Acquisizione dei dati	264
8.6.1. <i>Inserimento diretto</i>	264

8.6.2.	<i>Lettura di file di testo</i>	266
8.6.3.	<i>Lettura di file delimitati</i>	269
8.6.4.	<i>Lettura di file Excel</i>	270
8.6.5.	<i>Note</i>	273
9.	Analisi bivariata con il foglio elettronico	275
9.1.	Tabelle di contingenza o tabelle a doppia entrata	275
9.2.	Tabella pivot	278
9.2.1.	<i>Tabella pivot e suddivisione delle variabili in classi</i>	283
9.3.	Stima del rischio in Excel	285
9.3.1.	<i>Stima del rischio in Excel: tabella pivot</i>	285
9.3.2.	<i>Stima del rischio in Excel: riclassificazione e conteggio</i>	287
10.	Appendice	291
10.1.	Esempio di lettera di presentazione	291
10.2.	Esempio di questionario	293
10.3.	Tavole delle principali distribuzioni di probabilità	305
	Bibliografia	313

Prefazione

La scienza è fatta di dati come una casa è fatta di pietre.
Ma un ammasso di dati non è scienza più di quanto
un mucchio di pietre sia una vera casa.

Henri Poincarè

L'insieme degli argomenti trattati nel presente volume si propone come un essenziale ed agevole manuale operativo di *metodologia e tecnica della ricerca sociale*, affrontando con metodi assai pragmatici e concreti l'insieme degli argomenti di una attività di ricerca essenzialmente multidisciplinare.

È difficile coniugare gli aspetti metodologici caratteristici, diffusi, anche se condivisi con qualche distinguo, delle scienze della ricerca sociale, e gli strumenti che queste mettono a disposizione per la raccolta e l'elaborazione esplorativa e confermativa delle ipotesi formulate.

Gli argomenti attinenti le varie discipline vengono affrontati con un approccio semplificato, sufficiente a rendere autonomo il lettore, rinviando gli approfondimenti alla copiosa letteratura di settore, in parte citata in bibliografia.

Nel filo logico in cui la ricerca si esprime, si inseriscono ubiquitarie le nuove tecnologie, a prevalente derivazione informatica, sia hardware che software, con strumenti sempre più evoluti, che talvolta sopravvanzano lo sviluppo di un modello teorico di riferimento, introducendo implicazioni metodologiche non ancora comprese, o non comprese appieno.

Il ricercatore che ha maturato una qualche esperienza riesce a distinguere, a grandi linee, aspetti positivi e negativi delle nuove tecnologie, cosa che non riesce altrettanto bene a chi si avvicina solo marginalmente, o saltuariamente, a questi strumenti di ricerca sociale.

Il fervore riscontrato negli ultimi anni in questo settore della ricerca deve fare ora i conti sia con i rischi caratteristici di sempre, sia con gli errori introdotti da un uso superficiale dei nuovi prodotti, soprattutto quando questi promuovono indagini standardizzate, pronte per l'uso, con modelli o *template* preconfezionati.

L'esempio tipico del primo caso è l'indagine di opinione frettolosa, come ad esempio quella del giornalista che scende in strada ed intervista qualche passante, estrapolando dalle risposte ottenute l'opinione dell'intera popolazione in merito all'argomento del giorno.

Il secondo caso nasce paradossalmente dalla qualità e dalla quantità di strumenti software indirizzati specificatamente alla rilevazione ed alla analisi dei dati di indagine.

Proprio l'elevata qualità introdotta da parte di diversi applicativi, induce nel ricercatore una falsa sicurezza, portandolo a sviluppare l'attività di ricerca verso i paradigmi implicitamente inclusi nello strumento, senza porre accorta attenzione preliminare al disegno d'indagine ed alle diverse possibilità e metodi da utilizzare, con particolare riferimento agli obiettivi individuati nel momento della formulazione delle ipotesi iniziali.

Partendo da queste premesse, il presente lavoro vuole percorrere idealmente tutta la prassi oramai consolidata, almeno nelle sue linee più generali, di un percorso di ricerca sociale, con una particolare attenzione verso un suo specifico sottoinsieme, l'ambito socio – sanitario.

Idealmente il percorso si sviluppa attraverso le seguenti fasi:

- Scelta del problema e identificazione delle ipotesi d'indagine.
- Formulazione di un disegno della ricerca.
- Raccolta, codifica e validazione dei dati.
- Analisi esplorativa dei dati.
- Interpretazione dei risultati e analisi confermativa delle ipotesi iniziali.

È del tutto evidente che il compito che ci si prefigge è estremamente ampio, coinvolgendo nel suo percorso diverse discipline di derivazione sociologica, epidemiologica, informatica e statistica.

Per questo motivo gli argomenti vengono trattati ad un livello introduttivo, identificando come principale obiettivo il ruolo, non banale, di strumento multidisciplinare per la corretta implementazione ed attuazione di una indagine di ricerca.

I metodi di ricerca trattati sono essenzialmente quantitativi, ma non per questo vengono trascurati alcuni strumenti di tipo qualitativo che da sempre sono elementi fondanti della analisi sociologica, e che qui vengono visti come utili strumenti preliminari ed integrativi nella definizione degli obiettivi di indagine, dei temi da approfondire e delle modalità più idonee per farlo.

Tra questi vengono accennati:

- Intervista discorsiva.
- Osservazione e osservazione partecipata.
- Focus Group.
- Intervista con esperti, il metodo Delphi.

Grande risalto viene dato al momento della raccolta del dato, sia questo proveniente da fonti ufficiali o di tipo documentale, o costruito attraverso tecniche di intervista e di questionario.

Vengono inoltre definite caratteristiche, criticità ed opportunità delle varie modalità con cui viene somministrato il questionario:

- Indagine postale o email.
- Indagine faccia a faccia.
- Indagine via web.
- Indagine telefonica.
- Indagine con questionario autocompilato.

In tutte queste modalità di somministrazione del questionario esiste, più o meno diffusa, una versione *Computer Assisted*, ovvero con una integrazione di strumenti hardware e software che vanno dalla semplice opportunità di raccolta diretta del dato, alla somministrazione di questionari strutturati e personalizzati a seconda dell'utente-paziente-cliente a cui ci si rivolge.

Il mondo anglosassone, da sempre noto per la sua pragmaticità, non si smette neanche in questo settore, ed ha coniato per i diversi tipi di indagine i seguenti acronimi:

- CAPI *Computer Assisted Personal Interview.*
- PAPI *Paper Assisted Personal Interview.*
- CATI *Computer Assisted Telephone Interview.*
- CAWI *Computer Assisted Web Interview.*
- CASI *Computer Assisted Self Interview.*

In queste tipologie di indagine il questionario assume un ruolo determinante, non solo come strumento operativo per rilevare sul campione il valore o la modalità assunte dalle diverse variabili, ovvero le risposte registrate alle domande che vengono poste, ma come vero e proprio *framework* ove interpretare concettualmente le ipotesi iniziali del progetto di ricerca.

Nella redazione del questionario si deve tenere conto di una notevole quantità di elementi, come la sequenza delle domande, il modo di chiedere informazioni, la cautela nelle domande che causano imbarazzo, la congruenza semantica nelle possibilità di risposta, la unidimensionalità delle richieste, la comprensibilità del linguaggio utilizzato, l'aspetto grafico, ed altro ancora.

A tutti questi primari aspetti atti a tradurre nel concreto e nel migliore dei modi il disegno della ricerca, si aggiunge l'esigenza di particolari strutture di questionario che sono richieste a seconda degli strumenti di analisi statistica preventivamente ipotizzati.

In ambito epidemiologico, ad esempio, vi deve essere una domanda a risposta dicotomica, o dicotomizzabile, che permetta di distinguere tra casi e controlli, o tra esposti e non esposti, o tra malati e non malati, a seconda del tipo di analisi del rischio.

Oppure nella definizione di variabili latenti comuni a determinati tipologie di utenti-pazienti, nella ricerca di atteggiamenti attraverso la misura di opinioni espresse, si devono costruire delle domande multiple atte a rilevare delle scale, come ad esempio le scale di Likert, o le scale di Guttman o, delle scale aventi

una accezione probabilistica, come quelle identificate dal modello di Rasch e di Bradley-Terry.

La raccolta delle informazioni viene seguita da una fase di validazione e trasferimento su supporto informatico, operazione che richiedono una certa accuratezza ed una puntigliosa coerenza se si vuole pervenire ad una *Matrice di dati* di buona qualità e di facile fruibilità.

Fanno parte di questo momento l'analisi della coerenza dei dati raccolti, la ricerca di dati palesemente errati, ed il difficile passo ove vengono considerati i dati mancanti, compresa una loro eventuale imputazione.

L'accuratezza dei dati vorrebbe che, in presenza di una mancata risposta venga rimossa l'intera unità statistica, ma questo pesa inevitabilmente quando si hanno numerosità non molto elevate, o quando la rilevazione ha richiesto molte risorse, o quando si raccolgono dati su casi rari, come capita sovente in ambito sanitario in riferimento alle patologie poco frequenti nella popolazione.

Si preferisce talvolta imputare manualmente il dato mancante, sostituendolo con un valore neutro come potrebbe essere la mediana per dati almeno su scala ordinale, purchè i valori mancanti siano pochi punti percentuali rispetto al totale.

Nei casi di variabili qualitative si deve costruire una distribuzione di frequenza e valutare opportunamente caso per caso assegnando, ad esempio, la modalità più frequente, oppure alternando le due modalità più frequenti.

Il trasferimento su supporto informatico è quasi sempre necessario, anche per indagini svolte con interviste assistite dal computer, dato che si utilizzano sovente strutture di dati complesse, soprattutto in presenza di questionari strutturati.

La tecnologia ha fatto passi avanti notevoli nel settore dei sistemi informativi ma, in generale, i dati raccolti vengono trasferiti in un foglio elettronico per rilevazioni sino a circa 100 variabili e sino a circa 300 interviste, oltre si preferisce usufruire delle innegabili opportunità offerte da un gestore di base di dati, tipicamente una base di dati relazionale, utilizzando un *Relational Data Base Management System (RDBMS)*.

La fase di trasferimento delle informazioni è anche il momento più opportuno per definire il *dizionario dei dati (data dictionary)* o *libro dei codici (codebook)*, documento che raccoglie le definizioni dei nomi mnemonici dati alle variabili, come traduzione delle forme discorsive che queste assumono all'interno del questionario e le relative codifiche assegnate ad ognuna delle modalità possibili per la singola variabile.

Nel dizionario dei dati vengono annotate anche tutte le variazioni e le osservazioni di un qualche rilievo verificatesi nella somministrazione del questionario, e le ricodifiche o imputazioni di dati mancanti intervenute successivamente.

La successiva analisi dei dati può essere idealmente suddivisa in due momenti distinti, una prima fase comune di *analisi esplorativa*, valida per ogni tipo

di indagine, che utilizza strumenti statistici di base, ove si fanno classificazioni univariate del campione indagato, come distribuzioni per età, o territorio di provenienza, ed altro, utilizzando strumenti di statistica descrittiva, integrati da grafici ad istogramma, o di altro tipo.

Il secondo passo è la ricerca di associazione e di dipendenza tra due variabili, confrontando le variabili oggettive raccolte come età, genere, titolo di studio, ecc., con la o le variabili obiettivo dell'indagine che meglio identificano le ipotesi iniziali.

Lo strumento principale qui è la tabella bivariata, ove si confrontano le distribuzioni di frequenza assolute e relative percentuali per variabili categoriali o continue ridotte in classi, sino a giungere eventualmente alla sintesi massima delle tabelle tetracoriche.

Le ipotesi associative che sorgono vengono usualmente verificate, sia nella intensità puntuale, che nella significatività statistica attraverso il test chi – quadro e l'analisi del rischio utilizzando l'odds ratio quando risulti opportuno.

Anche in questo caso sono utili all'interpretazione i grafici ad istogrammi appaiati, oppure anche una verifica distributiva con boxplot sempre appaiati.

Il secondo momento di analisi dei dati, *l'analisi confermativa delle ipotesi*, si esplica attraverso una serie di test statistici, oppure attraverso la costruzione di veri e propri modelli statistici multivariati, i quali sono diversi da caso a caso, applicati a seconda del contesto, e sono qui solamente accennati.

L'aspetto operativo viene approfondito attraverso esempi pratici, tratti da lavori di indagine effettivamente svolti, secondo le diverse tipologie ed i diversi modelli proposti.

L'archiviazione dei dati e le elaborazioni più semplici vengono svolte utilizzando Calc di OpenOffice¹ o Excel di Microsoft Office, mentre le elaborazioni statistiche più complesse vengono attuate attraverso il software R².

¹ <http://www.openoffice.org>; <http://it.openoffice.org>

² <http://www.r-project.org>