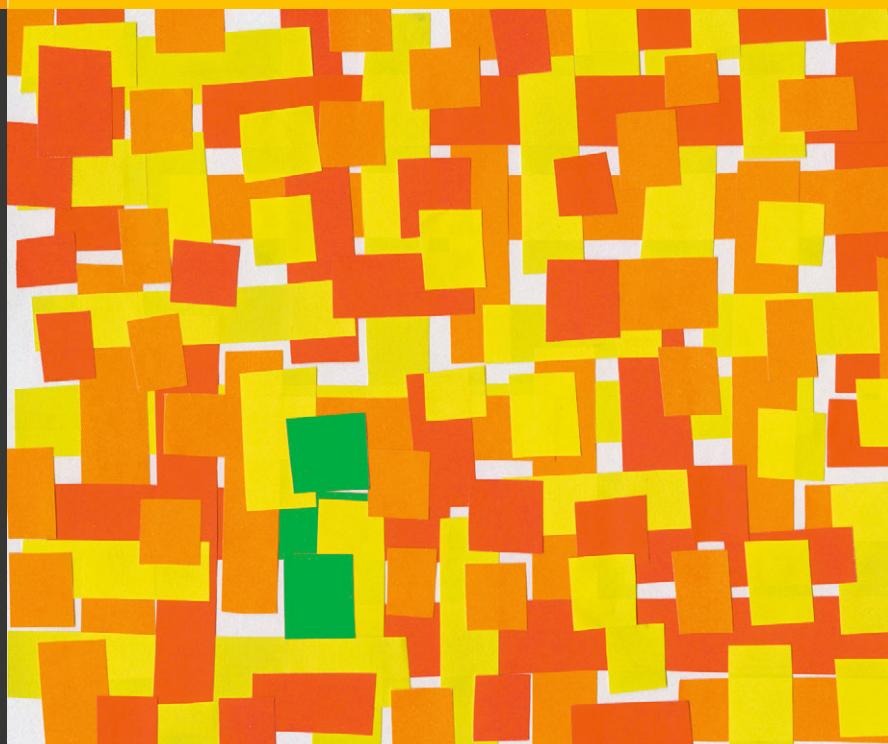


Giovanna Nicolini, Donata Marasini,  
Giorgio Eduardo Montanari, Monica Pratesi,  
Maria Giovanna Ranalli, Emilia Rocco

# Metodi di stima in presenza di errori non campionari



UNTEXT



Springer

## **Metodi di stima in presenza di errori non campionari**

Giovanna Nicolini • Donata Marasini •  
Giorgio Eduardo Montanari • Monica Pratesi •  
Maria Giovanna Ranalli • Emilia Rocco

# **Metodi di stima in presenza di errori non campionari**

Giovanna Nicolini  
Dipartimento di Economia,  
Management e Metodi Quantitativi  
Università degli Studi di Milano

Monica Pratesi  
Dipartimento di Statistica e Matematica  
applicata all'Economia  
Università degli Studi di Pisa

Donata Marasini  
Dipartimento di Economia,  
Metodi Quantitativi e Strategie di Impresa  
Università degli Studi di Milano-Bicocca

Maria Giovanna Ranalli  
Dipartimento di Economia,  
Finanza e Statistica  
Università degli Studi di Perugia

Giorgio Eduardo Montanari  
Dipartimento di Economia,  
Finanza e Statistica  
Università degli Studi di Perugia

Emilia Rocco  
Dipartimento di Statistica  
"G. Parenti"  
Università degli Studi di Firenze

UNITEXT – Collana di Statistica e Probabilità Applicata

ISSN 2240-2640

ISSN 2240-2659 (elettronico)

ISBN 978-88-470-2795-4

ISBN 978-88-470-2796-1 (eBook)

DOI 10.1007/978-88-470-2796-1

Springer Milan Heidelberg New York Dordrecht London

© Springer-Verlag Italia 2013

Quest'opera è protetta dalla legge sul diritto d'autore e la sua riproduzione anche parziale è ammessa esclusivamente nei limiti della stessa. Tutti i diritti, in particolare i diritti di traduzione, ristampa, riutilizzo di illustrazioni, recitazione, trasmissione radiotelevisiva, riproduzione su microfilm o altri supporti, inclusione in database o software, adattamento elettronico, o con altri mezzi oggi conosciuti o sviluppati in futuro, rimangono riservati. Sono esclusi brevi stralci utilizzati a fini didattici e materiale fornito ad uso esclusivo dell'acquirente dell'opera per utilizzazione su computer. I permessi di riproduzione devono essere autorizzati da Springer e possono essere richiesti attraverso RightsLink (Copyright Clearance Center). La violazione delle norme comporta le sanzioni previste dalla legge.

Le fotocopie per uso personale possono essere effettuate nei limiti del 15% di ciascun volume dietro pagamento alla SIAE del compenso previsto dalla legge, mentre quelle per finalità di carattere professionale, economico o commerciale possono essere effettuate a seguito di specifica autorizzazione rilasciata da CLEARedi, Centro Licenze e Autorizzazioni per le Riproduzioni Editoriali, e-mail autorizzazioni@clearedi.org e sito web [www.clearedi.org](http://www.clearedi.org).

L'utilizzo in questa pubblicazione di denominazioni generiche, nomi commerciali, marchi registrati, ecc., anche se non specificatamente identificati, non implica che tali denominazioni o marchi non siano protetti dalle relative leggi e regolamenti.

Le informazioni contenute nel libro sono da ritenersi veritiere ed esatte al momento della pubblicazione; tuttavia, gli autori, i curatori e l'editore declinano ogni responsabilità legale per qualsiasi involontario errore od omissione. L'editore non può quindi fornire alcuna garanzia circa i contenuti dell'opera.

9 8 7 6 5 4 3 2 1

*Layout copertina:* Beatrice E, Milano

Impaginazione: PTP-Berlin, Protogal T<sub>E</sub>X-Production GmbH, Germany ([www.ptp-berlin.eu](http://www.ptp-berlin.eu))

Springer-Verlag Italia S.r.l., Via Decembrio 28, I-20137 Milano

Springer-Verlag fa parte di Springer Science+Business Media ([www.springer.com](http://www.springer.com))

---

## Prefazione

Questo volume nasce dalla collaborazione di sei ricercatori, tutti di estrazione accademica e con un comune interesse nella metodologia delle indagini campionarie. Gli autori si sono trovati ad organizzare nel 2007 un corso dal titolo “Teoria e pratica delle indagini campionarie: approccio probabilistico ed errori non campionari” che si è svolto all’Università di Milano Bicocca nell’ambito dei corsi di alta formazione patrocinati dalla Società Italiana di Statistica. La teoria dei campioni rappresenta una parte della Statistica molto articolata e ogni autore ha apportato al corso un suo specifico contributo che, successivamente, si è trasformato in un argomento trattato nel presente libro.

La finalità di questo volume è quella di dare ai ricercatori gli strumenti necessari per correggere gli errori non campionari che inevitabilmente sorgono nelle diverse tipologie di indagine, mettendo in relazione tali errori anche con i più moderni metodi di condurre le indagini. Per come è stato impostato e per il rigore scientifico impiegato – che colloca nelle Dimostrazioni le parti metodologiche più sofisticate – il target di riferimento può essere sia accademico (corsi in lauree specialistiche, dottorati, master) sia non accademico (istituti di ricerca in generale). Il volume riporta argomenti che si trovano in letteratura, tuttavia il suo valore aggiunto si ritiene sia quello di averli unificati e di averli inseriti in un contesto moderno di indagine.

I primi due capitoli, pur non essendo dedicati agli errori non campionari, costituiscono un utile link con i metodi più tradizionali di campionamento e stima e unificano simboli e terminologie. Il Cap. 3 presenta una trattazione sistematica dell’utilizzo delle variabili ausiliarie nell’inferenza descrittiva su popolazioni finite; questo argomento non è facilmente reperibile in letteratura e costituisce la base necessaria per lo sviluppo dei capitoli successivi. Inoltre in questo capitolo si fa un accenno all’approccio basato sul modello, anch’esso indispensabile, per comprendere appieno i metodi proposti, e si analizzano gli effetti degli errori di misura, in modo da evidenziare la necessità della loro prevenzione, rinviando però ai testi specializzati i metodi per il loro trattamento.

I successivi tre capitoli sono dedicati al trattamento degli errori di copertura e di mancata osservazione e dei loro effetti sui risultati di una indagine statistica. Il Cap. 4 si occupa in particolare dei problemi connessi alla disponibilità di liste delle unità da campionare e alla loro capacità di rappresentare correttamente la popolazione indagata. Si analizzano le circostanze che portano ad errori e le loro ripercussioni sulla qualità delle stime ottenute e vengono suggeriti i metodi più utilizzati per prevenire e/o correggere gli errori di copertura, anche nei casi, sempre più frequenti, delle indagini via web. I Capp. 5 e 6 trattano il caso della mancata osservazione delle unità designate a far parte del campione, cioè la non risposta. Il Cap. 5 analizza gli effetti e presenta i metodi di correzione della non risposta totale basati principalmente sulla riponderazione dei dati e, tra questi, particolare attenzione è riservata ai metodi di calibrazione. Il Cap. 6 si occupa invece della non risposta parziale e delle tecniche di imputazione correntemente utilizzate per farvi fronte, senza tralasciare argomenti più avanzati frutto della ricerca più recente sul tema.

Al fine di agevolare la comprensione di quanto illustrato, in ogni capitolo sono stati inseriti diversi esempi, tratti a volte da indagini reali, per declinare nel caso considerato le metodologie descritte.

Tutti i capitoli sono stati ampiamente discussi dagli autori che hanno condiviso l'impostazione metodologica, tuttavia i Capp. 1 e 2 sono stati scritti da Donata Marasini e Giovanna Nicolini, il Cap. 3 da Giorgio E. Montanari, Monica Pratesi e Maria Giovanna Ranalli, il Cap. 4 da Giovanna Nicolini e Monica Pratesi, il Cap. 5 da Giorgio E. Montanari e Maria Giovanna Ranalli e, infine, il Cap. 6 da Emilia Rocco.

Milano, settembre 2012

Giovanna Nicolini  
Donata Marasini  
Giorgio Eduardo Montanari  
Monica Pratesi  
Maria Giovanna Ranalli  
Emilia Rocco

---

# Indice

<b>1</b>	<b>Introduzione al campionamento da popolazioni finite</b>	<b>1</b>
1.1	Introduzione	1
1.2	Indagine campionaria	1
1.2.1	Definizioni preliminari	2
1.2.2	Campione casuale, probabilistico, rappresentativo	5
1.2.3	Fasi di un'indagine campionaria	6
1.2.4	Gli errori di un'indagine campionaria	8
1.3	Impiego di variabili ausiliarie nei piani di campionamento probabilistici	10
1.3.1	Campionamento a due fasi	12
1.4	Campionamenti non probabilistici	13
	Bibliografia	14
<b>2</b>	<b>I campionamenti probabilistici</b>	<b>17</b>
2.1	Introduzione	17
2.2	Piano di campionamento	17
2.3	Le probabilità di inclusione	19
2.4	La variabile indicatrice e la frequenza attesa di inclusione	20
2.5	Alcuni piani di campionamento	22
2.5.1	Piano di campionamento casuale semplice senza reinserimento e con reinserimento	22
2.5.2	Piano di campionamento sistematico ( <i>sm</i> )	23
2.5.3	Piano di campionamento di Poisson ( <i>Po</i> )	24
2.5.4	Piano di campionamento con reinserimento proporzionale alla misura di ampiezza ( <i>pps</i> )	25
2.5.5	Piani di campionamento senza reinserimento proporzionali alla misura di ampiezza ( <i>pps</i> ) per $n = 2$	25
2.5.6	Piani di campionamento $\pi ps$ per $n > 2$	26
2.5.7	Piano di campionamento stratificato ( <i>st</i> )	28
2.5.8	Piano di campionamento a grappoli ( <i>gr</i> )	29
2.5.9	Piano di campionamento a due stadi ( <i>ds</i> )	30
2.5.10	I piani di campionamento complessi	30

2.6	Stime e stimatori .....	32
2.6.1	Lo stimatore di Horvitz-Thompson .....	32
2.6.2	Le varianze .....	34
2.6.3	Lo stimatore di Hansen-Hurwitz .....	37
2.7	Effetto del disegno ed efficienza .....	38
2.8	Stimatori di parametri funzioni di totali .....	39
2.9	La stima di un quantile .....	41
	Dimostrazioni .....	42
2.10	Le probabilità di inclusione .....	42
2.11	La correttezza dello stimatore della varianza dello stimatore di Horvitz-Thompson .....	43
2.12	La correttezza e la varianza dello stimatore di Horvitz-Thompson nel campionamento $ds$ .....	43
2.13	Il metodo della linearizzazione .....	44
	Bibliografia .....	45
<b>3</b>	<b>L'impiego delle variabili ausiliarie per la costruzione degli stimatori</b> .....	47
3.1	Introduzione .....	47
3.2	La costruzione degli stimatori nell'approccio <i>design-based</i> ....	49
3.2.1	Lo stimatore per quoziente e le sue proprietà .....	50
3.2.2	Lo stimatore post-stratificato e le sue proprietà .....	54
3.2.3	Lo stimatore per differenza e le sue proprietà .....	58
3.2.4	La stima per regressione generalizzata .....	65
3.2.5	La ponderazione dei dati nelle indagini campionarie ....	72
3.2.6	La stima a ponderazione vincolata o calibrazione .....	73
3.2.7	La stima dei parametri nei domini di studio .....	81
3.3	La costruzione degli stimatori nell'approccio <i>model-based</i> ....	85
3.3.1	Obiettivo dell'inferenza nell'approccio basato sul modello .....	85
3.3.2	Stima dei parametri descrittivi .....	86
3.3.3	Stima dei parametri di superpopolazione .....	89
3.3.4	Interpretazione del concetto di superpopolazione .....	91
3.3.5	Predittori o stimatori? .....	93
3.4	Inferenza da modello ed errori di misura .....	94
3.4.1	Gli effetti degli errori di misura .....	95
3.4.2	Effetto degli errori di misura sul valore atteso degli stimatori .....	96
3.4.3	Effetti degli errori di misura sulla varianza delle stime ..	97
3.4.4	Effetti degli errori di misura sugli stimatori della varianza .....	99
	Dimostrazioni .....	101
3.5	La procedura per la determinazione dei pesi calibrati .....	101
	Bibliografia .....	103



<b>4</b>	<b>Metodi per correggere gli errori di copertura</b>	105
4.1	Introduzione	105
4.2	Copertura e modo di indagine	106
4.2.1	Interviste faccia a faccia e liste di grappoli di unità della popolazione oggetto d'indagine	106
4.2.2	Interviste telefoniche ed elenchi di numeri di telefono	107
4.2.3	Indagini postali, elenchi anagrafici ed elenchi elettorali	109
4.2.4	Indagini Web, liste di indirizzi di posta elettronica e di utenti Internet	110
4.3	Errori nella popolazione frame	111
4.3.1	Tipologie di frame	112
4.3.2	Tipologia di errori nel frame	114
4.3.3	Conseguenza degli errori nel frame	116
4.4	Metodi per contenere gli errori di copertura	121
4.4.1	Frame ad hoc	121
4.4.2	Frame da fonti amministrative	122
4.4.3	Scelta tra più frame	123
4.4.4	Frame costruiti durante l'indagine	123
4.5	Metodi di indagini con frame imperfetti	125
4.5.1	Campionamenti con molteplicità	125
4.5.2	Uso congiunto di più frame	130
	Bibliografia	139
<b>5</b>	<b>Metodi inferenziali in presenza di mancate risposte totali</b>	143
5.1	Introduzione	143
5.2	Interventi preventivi a livello del disegno dell'indagine	146
5.3	Tecniche di correzione della distorsione da mancata risposta in fase di stima	148
5.3.1	Classi di aggiustamento per mancate risposte	149
5.3.2	Approccio in due fasi	153
5.3.3	Stima delle probabilità di risposta attraverso il modello logistico	157
5.3.4	Stima per regressione	159
5.3.5	Stimatori a ponderazione vincolata in presenza di mancate risposte	164
5.4	Schemi di mancata risposta e selezione delle variabili ausiliarie	169
	Bibliografia	172
<b>6</b>	<b>Metodi inferenziali in presenza di mancate risposte parziali</b>	173
6.1	Introduzione	173
6.2	Tipologie e cause della mancata risposta parziale	174
6.3	Azioni per prevenire/ridurre le mancate risposte parziali	175
6.4	Metodi per trattare i dati incompleti	177
6.4.1	Metodi basati sulle sole unità rispondenti	179
6.4.2	Metodi di riponderazione	180

6.4.3	Metodi di imputazione .....	182
6.5	Imputazione multipla .....	199
6.6	La stima della varianza in presenza di valori imputati .....	202
6.6.1	Metodi numerici di stima della varianza .....	207
6.7	Note di approfondimento sull'imputazione multipla .....	209
6.7.1	Giustificazione teorica dell'imputazione multipla .....	209
6.7.2	Proprietà dell'imputazione multipla nell'ambito dell'inferenza randomizzata .....	212
6.7.3	Metodi Bayesiani iterativi per realizzare imputazioni multiple .....	215
6.7.4	Imputazione multipla per meccanismi di risposta non ignorabili e <i>sensitivity analysis</i> .....	216
	Bibliografia .....	218
	<b>Indice analitico</b> .....	221

---

# Simboli

---

Popolazione	
$U$	Popolazione obiettivo (o target)
$N$	Numero di unità statistiche o elementari che formano la popolazione finita
$F$	Lista della popolazione (o popolazione frame)
$M$	Numero delle unità di rilevazione che formano la lista della popolazione
$N_h$	Numero di unità statistiche o elementari nella popolazione che formano lo strato $h$ -esimo
$h = 1, \dots, H$	Indice di strato
$M_i$	Numero di unità statistiche o elementari nella popolazione che formano il grappolo $i$ -esimo
$i = 1, \dots, M$	Indice di grappolo
$N_l$	Numero di unità statistiche o elementari nella popolazione che formano il post-strato $l$ -esimo
$l = 1, \dots, L$	Indice di post-strato
$N_D$	Numero di unità statistiche o elementari nella popolazione che formano il dominio $D$ -esimo
$N_F$	Numero delle unità statistiche della popolazione target contenute nel frame
$N_{F^*}$	Numero delle unità statistiche non appartenenti alla target ma presenti nel frame
$N_M$	Numero delle unità statistiche nella popolazione target mancanti nel frame
$W_M$	Tasso di sotto-copertura
$W_{F^*}$	Tasso di sovra-copertura
$N_R$	Numero di unità statistiche o elementari nella sotto-popolazione dei rispondenti
$N_{NR}$	Numero di unità statistiche o elementari nella popolazione che formano la sottopopolazione dei non rispondenti

**Variabili**

$y$	Variabile oggetto di studio
$y_j$	Valore della variabile $y$ riferito alla $j$ -esima unità della popolazione
$x$	Variabile ausiliaria
$x_j$	Valore della variabile $x$ riferito alla $j$ -esima unità della popolazione
$I_j(s)$	Variabile indicatrice di appartenenza ad $s$
$\mathbf{x} = (x_1, x_2, \dots, x_P)$	Vettore di variabili ausiliarie
$\mathbf{x}_j$	Vettore dei valori di $\mathbf{x}$ nella $j$ -esima unità
$\mathbf{y} = (y_1, y_2, \dots, y_N)$	Vettore degli $N$ valori di $y$ nella popolazione
$\mathbf{Y} = (Y_1, Y_2, \dots, Y_N)$	Vettore delle $N$ variabili casuali $Y_j$ associate alle unità della popolazione
$\phi$	Variabile molteplicità
$\phi_j$	Valore della molteplicità dell'unità statistica $j$ -esima
$N_\phi$	Frequenza assoluta della molteplicità con valore $\phi$
$Y_\phi$	Totale della variabile $y$ sulle unità che hanno una molteplicità pari a $\phi$
$R_j$	Variabile indicatrice di appartenenza all'insieme dei rispondenti $r$
$\mathbf{y}_\bullet = \{y_{\bullet j} : j \in r\}$	Insieme completato dei dati mediante un qualsiasi metodo di imputazione

**Parametri**

$T$	Parametro della popolazione
$Y$	Totale nella popolazione della variabile $y$
$\bar{Y}$	Valore medio nella popolazione di $y$
$A$	Frequenza assoluta nella popolazione
$P$	Frequenza relativa o proporzione nella popolazione
$C_y$	Coefficiente di variazione nella popolazione di $y$
$Q_q$	Quantile di ordine $q$
$\Phi(y)$	Funzione di ripartizione nella popolazione della variabile $y$
$X$	Totale della variabile ausiliaria $x$
$\bar{X}$	Valore medio nella popolazione di $x$
$\mathbf{X} = (X_1, X_2, \dots, X_P)$	Vettore dei totali nella popolazione di $x$
$S_y^2$	Varianza nella popolazione di $y$
$S_x^2$	Varianza nella popolazione di $x$
$S_{yx}$	Covarianza tra le variabili $y$ e $x$
$\beta_{yx}$	Coefficiente di regressione
$\rho_{yx}$	Coefficiente di correlazione lineare
$R$	Rapporto tra totali o tra medie
$\theta$	Vettore dei parametri di superpopolazione
$f(\cdot)$	Funzione di densità di una variabile casuale
$\varepsilon \quad \sigma^2$	Componente casuale del modello di superpopolazione e sua varianza
$\mu_j \quad \sigma_j^2 \quad \sigma_{jj'}$	Rispettivamente valore atteso e varianza di $Y_j$ e covarianza tra $Y_j$ e $Y_{j'}$ secondo il modello di superpopolazione

Campione	
$s$	Generico campione
$n$	Dimensione fissa di $s$
$n(s)$	Dimensione variabile di $s$
$f = n/N$	Frazione di sondaggio
$S_0$	Spazio dei campioni ordinati, con e senza reinserimento, di ampiezza variabile
$S$	Spazio dei campioni non ordinati e senza reinserimento di ampiezza variabile
$p(s)$	Piano di campionamento (pdc)
$p_j$	Probabilità di selezione della $j$ -esima unità della popolazione
$\pi_j$	Probabilità di inclusione del primo ordine dell'unità $j$ -esima
$\pi_{j,j'}$	Probabilità di inclusione del secondo ordine delle unità $j$ e $j'$
$F_j(s)$	Frequenza attesa di inclusione
$\psi_j$	Numero medio di volte in cui l'etichetta $j$ è compresa in $s$
$a_j = 1/\pi_j$	Peso base di riporto all'universo
$a_{j,j'} = 1/\pi_{j,j'}$	Reciproco della probabilità di inclusione del secondo ordine
$\pi_k$	Probabilità di inclusione del primo ordine del network $k$
$\pi_{kl}$	Probabilità di inclusione del secondo ordine dei network $k$ e $l$
$r$	Campione dei rispondenti
$n_r$	Dimensione del campione dei rispondenti
$q(r s)$	Processo di selezione dei rispondenti dal campione $s$
$\vartheta_j$	Probabilità di risposta dell'unità $j$ -esima
$r_c$	Insieme delle unità completamente osservate
$r_y$	Insieme delle unità che hanno fornito un valore per la variabile di interesse $y$
$n_{r_y}$	Dimensione insieme delle unità che hanno fornito un valore per la variabile di interesse $y$
Piani di Campionamento	
$srs$	pdc casuale semplice senza reinserimento
$srswr$	pdc casuale semplice con reinserimento
$sm$	pdc sistematico
$Po$	pdc di Poisson
$pps$	pdc con reinserimento proporzionale alla misura di ampiezza
$\pi ps$	pdc senza reinserimento proporzionale alla misura di ampiezza
$st$	pdc stratificato
$gr$	pdc a grappoli
$ds$	pdc a due stadi

---

**Stimatori**


---

$\hat{T}$	Stimatore generico del parametro $T$
$\hat{Y}_\pi$ ,	Stimatore o stima di Horvitz-Thompson (HT) del totale $Y$
$\hat{X}_\pi$	Stimatore o stima di HT del totale $X$
$\hat{R}_\pi$	Stimatore o stima di HT del rapporto $R$
$\hat{\bar{Y}}_\pi, \hat{\bar{X}}_\pi$	Stimatore o stima di HT di $\bar{Y}$ e di $\bar{X}$
$\hat{Y}, \hat{X}, \hat{A}$	Stimatore di HT di $Y$ , di $X$ e di $A$ per un <i>srs</i>
$\bar{y}, \bar{x}$	Media campionaria di $y$ e di $x$ .
$\hat{Y}_{HH}$	Stimatore o stima di Hansen-Hurwitz
$\hat{Y}_{PS}, \hat{Y}_{PSS}$	Stimatore post-stratificato in presenza di mancate risposte
$\hat{Y}_{\pi st}$	Stimatore o stima di HT del totale $Y$ per un <i>st</i>
$\hat{Y}_{\pi gr}$	Stimatore o stima di HT del totale $Y$ per un <i>gr</i>
$\hat{Y}_{\pi ds}$	Stimatore o stima di HT del totale $Y$ per un <i>ds</i>
$\hat{Y}_{2f}$	Stimatore nel campionamento a due fasi
$\hat{Q}_q$	Stimatore o stima di un quantile
$\hat{Y}_q$ ,	Stimatore o stima per quoziente del totale $Y$
$\hat{Y}_{ps}$	Stimatore o stima post-stratificata del totale $Y$
$\hat{Y}_d$	Stimatore o stima per differenza del totale $Y$
$\hat{Y}_{reg}$	Stimatore o stima per regressione generalizzata del totale $Y$
$\hat{Y}_{cal}$	Stimatore o stima di calibrazione del totale $Y$
$\hat{Y}_\xi$	Stimatore basato sul modello del totale $Y$
$\hat{Y}_{Ha}$	Stimatore di Hartley
$\hat{Y}_M$	Stimatore del totale in presenza di molteplicità
$\hat{Y}_H$	Stimatore di Hajek
$\hat{T}_\bullet$	Stimatore del parametro $T$ applicato all'insieme completo dei dati

---

**Operatori**


---

$B(\cdot)$	Distorsione rispetto al piano di campionamento
$E(\cdot)$	Valore atteso rispetto al piano di campionamento
$V(\cdot)$	Varianza rispetto al piano di campionamento
$C(\cdot, \cdot)$	Covarianza rispetto al piano di campionamento
$v(\cdot)$	Stima della varianza
$c(\cdot, \cdot)$	Stima della covarianza
$Deff$	Effetto del disegno
$Eff$	Efficienza relativa
MSE	Errore quadratico medio
$\xi$	Modello di superpopolazione
$E_\xi(\cdot)$	Valore atteso rispetto al modello di superpopolazione
$V_\xi(\cdot)$	Varianza rispetto al modello di superpopolazione
$V_{p\xi}(\cdot)$	Varianza congiunta rispetto al piano di campionamento ed al modello di superpopolazione
$C_\xi(\cdot, \cdot)$	Covarianza rispetto al modello di superpopolazione
$E_m(\cdot)$	Valore atteso rispetto al modello di misura
$V_m(\cdot)$	Varianza rispetto al modello di misura

$V_{pm}(\cdot)$	Varianza congiunta rispetto al piano di campionamento e al modello di misura
$E_q(\cdot s)$	Valore atteso rispetto al processo di risposta, condizionato al campione estratto
$V_q(\cdot s)$	Varianza rispetto al processo di risposta, condizionato al campione estratto
$V_{pq}(\cdot)$	Varianza congiunta rispetto al piano di campionamento ed al processo di risposta
$v_q(\cdot s)$	Stima della varianza rispetto al processo di risposta, condizionato al campione estratto
$E_{MI}$	Valore atteso rispetto a un modello di imputazione
$V_{MI}$	Varianza rispetto a un modello di imputazione
$v_{JK}(\cdot)$	Stima Jackknife della varianza
$v_{bo}(\cdot)$	Stima bootstrap della varianza

---

# Introduzione al campionamento da popolazioni finite

## 1.1 Introduzione

Lo scopo introduttivo di questo capitolo è quello di richiamare i concetti fondamentali del campionamento da popolazioni finite. Dopo alcune definizioni preliminari e una breve disquisizione sulla terminologia in uso, vengono ricordate le fasi di una indagine statistica e individuati gli errori che in ciascuna di esse possono aver origine. Come si vedrà nei successivi capitoli di questo volume, dedicati proprio alle diverse tipologie di errore, l'impiego delle variabili ausiliarie per la individuazione e la correzione degli errori è fondamentale. Tuttavia, come è noto, ad esse si può ricorrere anche per la costruzione del piano di campionamento. Questo non implica la correzione di alcun errore, ma una maggiore precisione ed aderenza alla realtà che si vuole indagare e il ricorso alle variabili ausiliarie in questo caso può essere inteso come un metodo preventivo degli errori. Non sempre è possibile disporre delle variabili ausiliarie per i diversi impieghi, può allora essere necessario costruire ad hoc un data-set di tali variabili, per esempio con il campionamento a due fasi, che viene brevemente richiamato. Infine, vengono ricordati alcuni metodi di indagine che non seguono i canoni tradizionali del campionamento da popolazioni finite – sono i così detti *campionamenti non probabilistici* – alcuni molto in uso nelle indagini in ambito sociale. Anche se per essi non è possibile conoscere l'errore campionario, tuttavia, in questi casi, si può ricorrere a metodi inferenziali diversi da quello tradizionale usato nel campionamento da popolazioni finite, noto come *design-based*.

## 1.2 Indagine campionaria

Si supponga di voler studiare una variabile  $y$  su una popolazione finita formata da  $N$  unità statistiche, con lo scopo di conoscerne il totale, il suo valor medio o una qualsiasi altra funzione dei dati che ne definisce un *parametro*  $T$ . Se si effettua un'*indagine censuaria* tutte le unità della popolazione vengono



osservate e, nell'ipotesi che non si verifichi alcun tipo di errore nella rilevazione dei dati, si determina il reale valore del parametro di interesse  $T$ . Se invece si effettua una *indagine campionaria* solo una parte delle unità della popolazione viene osservata; ne consegue che, anche in assenza di errori di rilevazione, il valore reale del parametro è ignoto, tuttavia può essere stimato attraverso i valori delle unità che costituiscono il campione.

La non conoscenza del reale valore del parametro comporta la necessità di una sua stima che deve essere la “migliore”; questa esigenza pone il problema dell'attendibilità della stima e della scelta delle tecniche di campionamento. La presenza di anomalie nella costruzione del piano di campionamento e/o della stima allontanano quest'ultima dal valore reale del parametro. Pertanto, in un'indagine concepita secondo la tradizione classica (basata sul disegno), la conoscenza e la identificabilità delle unità della popolazione di riferimento sono fondamentali, come fondamentali sono alcuni presupposti logici presentati nel seguito.

### 1.2.1 Definizioni preliminari

Il campionamento da popolazione finita si basa su presupposti logici che sono ormai noti in letteratura (Hansen et al. 1953; Cochran 1977; Kish 1965, 1987; Särndal et al. 1992; Frosini et al. 2011); tuttavia, qui di seguito, se ne vogliono richiamare alcuni ai quali si farà riferimento nel corso del presente volume:

- *Popolazione obiettivo* o *popolazione target*. È la popolazione formata da un numero finito  $N$  di *unità statistiche* identificabili, su cui si analizzano la variabile o le variabili oggetto dell'indagine. Se ad ogni unità si associa un numero intero da 1 a  $N$ , tale popolazione viene indicata con  $U = \{1, \dots, j, \dots, N\}$ .
- *Lista della popolazione* o *popolazione frame*. È l'elenco delle  $M$  *unità di rilevazione* e viene indicato con  $F = \{F_1, \dots, F_i, \dots, F_M\}$ . L'unità di rilevazione può coincidere con l'unità statistica – in tal caso  $N = M$  – ovvero è un grappolo di unità statistiche e pertanto ogni unità di rilevazione  $F_i$  è una subpopolazione composta da  $M_i$  unità statistiche, il cui totale è  $\sum_{i=1}^M M_i = N$ . Ad esempio, nel primo caso la popolazione target è formata dagli individui residenti in un comune e la popolazione frame ne è l'elenco anagrafico; nel secondo caso la popolazione target è ancora formata dagli individui residenti e quella frame è l'elenco dalle famiglie residenti in quel comune. Rientra nel secondo caso il così detto *frame areale*, le cui unità di rilevazione sono le *areole* che rappresentano parti ben definite sul territorio di riferimento quali, ad esempio, i quartieri di una città, le sezioni di censimento o anche parti di territorio definite ad hoc. Ogni areola contiene un numero differente di unità statistiche che possono essere, a seconda dei casi, famiglie, individui, aziende agricole, unità commerciali, ecc. (per una analisi più approfondita delle relazioni tra popolazioni frame e target si veda il Cap. 4).

- *Parametri della popolazione.* Indicato con  $y_j$  ( $j = 1, \dots, N$ ) il valore della variabile in esame  $y$  riferito alla  $j$ -esima unità della popolazione, si definisce parametro una qualsiasi funzione degli  $N$  valori:  $T = f(y_1, \dots, y_N)$ . Numerosi sono i parametri di una variabile, nel seguito vengono riportati quelli descrittivi più comunemente usati. Nel caso di una variabile quantitativa, il parametro di maggior interesse è il *totale*

$$Y = \sum_{j=1}^N y_j, \quad (1.1)$$

o anche il *valor medio* che è dato dal rapporto tra totale e dimensione della popolazione:

$$\bar{Y} = \frac{Y}{N}. \quad (1.2)$$

Se invece la variabile è qualitativa e dicotoma, i due parametri più utilizzati sono ancora il totale e il valor medio che, tuttavia, in questo contesto, equivalgono rispettivamente ad una *frequenza assoluta*  $A$  e ad una frequenza relativa o *proporzione*  $P$ . Infatti, per un fenomeno qualitativo, indicata con  $C$  la modalità di interesse e con  $\bar{C}$  la modalità associata a tutto ciò che ne è complementare, se l'unità  $j$  possiede  $C$  si conviene di porre  $y_j = 1$ , mentre se possiede  $\bar{C}$ , si pone  $y_j = 0$ . Ne consegue che il numero di unità della popolazione con la modalità  $C$  è pari a

$$A = \sum_{j=1}^N y_j = Y \quad (1.3)$$

e rappresenta una frequenza assoluta. Dividendo la (1.3) per la dimensione della popolazione si ottiene una frequenza relativa, cioè la proporzione delle unità della popolazione con la caratteristica di interesse

$$P = \frac{A}{N}. \quad (1.4)$$

Parametri che identificano posizioni particolari sulla distribuzione della variabile sono i *quantili*. Indicata con  $\Phi(y)$  la funzione di ripartizione della variabile  $y$ , sempre che la medesima possa essere ordinata, il *quantile di ordine*  $q$  ( $0 < q < 1$ ), in termini operativi, si può definire nel modo seguente:

$$Q_q = \inf\{y : \Phi(y) \geq q\}. \quad (1.5)$$

Un particolare quantile molto utilizzato per la sintesi di distribuzioni asimmetriche è la *mediana* che si ha per  $q = 0,5$ . Un altro parametro di interesse per una variabile quantitativa è la *varianza*, che esprime la variabilità della

variabile  $y$  nella popolazione, data in questo contesto da (Cochran 1997)<sup>1</sup>

$$S_y^2 = \frac{1}{N-1} \sum_{j=1}^N (y_j - \bar{Y})^2 \quad (1.6)$$

e la sua radice quadrata  $S_y$ , nota come scarto quadratico medio. Un parametro molto usato per confrontare la variabilità è il *coefficiente di variazione* dato dal rapporto

$$C_y = \frac{S_y}{\bar{Y}}. \quad (1.7)$$

Oltre ai parametri funzioni di una sola variabile potrebbero interessare anche quelli che mettono in evidenza la relazione tra la variabile  $y$  e una variabile  $x$  presente sulle unità della popolazione, come il *rapporto* tra totali o tra medie

$$R = \frac{Y}{X} = \frac{\bar{Y}}{\bar{X}}, \quad (1.8)$$

con  $\sum_{j=1}^N x_j = X$  totale della variabile  $x$  ; il *coefficiente di regressione*

$$\beta_{yx} = \frac{S_{yx}}{S_x^2}, \quad (1.9)$$

dove  $S_x^2$  è la varianza della variabile  $x$  e  $S_{yx}$  è la covarianza tra le due variabili

$$S_{yx} = \frac{1}{N-1} \sum_{j=1}^N (y_j - \bar{Y})(x_j - \bar{X}) \quad (1.10)$$

e il *coefficiente di correlazione lineare*

$$\rho_{yx} = \frac{S_{yx}}{S_y S_x}. \quad (1.11)$$

Vale la pena notare che i parametri descritti possono essere interpretati come funzioni di uno o più totali relativi alla sola variabile  $y$  o alla coppia  $(y, x)$ . Infatti, è immediato osservare che i parametri (1.2) e (1.4) sono funzioni del solo totale  $Y$ , mentre è facile verificare che il parametro (1.6) è funzione dei totali  $\sum_{j=1}^N y_j^2$  e  $Y$  potendosi scrivere

$$S_y^2 = \frac{1}{N-1} \sum_{j=1}^N y_j^2 - \frac{1}{N(N-1)} Y^2. \quad (1.12)$$

I parametri (1.8), (1.9), (1.10) e (1.11) chiamano invece in causa totali relativi anche alla variabile  $x$  e alla coppia  $(y, x)$ .

---

<sup>1</sup> Conviene notare che nel presente contesto i parametri (1.6) e (1.10) sono rapportati alla quantità  $N-1$ , anziché a  $N$  come è naturale, esclusivamente per motivi di semplicità quando si passa alla stima dei medesimi.

### 1.2.2 Campione casuale, probabilistico, rappresentativo

*Campione casuale, probabilistico, rappresentativo*: sono tre concetti fondamentali nel campionamento da popolazioni finite che non necessariamente devono coesistere, avendo ciascuno la sua specificità.

La casualità riguarda il modo in cui sono selezionate le unità della popolazione destinate a formare il campione; pertanto la selezione deve avvenire seguendo una qualsiasi procedura riconducibile ad un esperimento aleatorio. Se la dimensione  $N$  della popolazione è nota, l'estrazione casuale è assimilabile all'estrazione di palline da un'urna che contiene appunto  $N$  palline; se invece la dimensione della popolazione non è nota, come ad esempio possono essere i clienti di un grosso emporio di abbigliamento, è sempre possibile pensare ad un meccanismo di casualità che selezioni alcune unità; ad esempio si può pensare di intervistare un cliente che entra nell'emporio ogni  $m$  ingressi.

Un campione è probabilistico quando la probabilità che ciascuna unità della popolazione ha di farne parte, in base al meccanismo di selezione, è positiva e calcolabile, secondo la definizione fornita nel prossimo capitolo. In questo senso un campione probabilistico è anche casuale.

La rappresentatività di un campione, invece, è un concetto più articolato al quale si possono ricondurre più definizioni. Kruskal e Mosteller (1979) hanno individuato sei definizioni di rappresentatività presenti nella letteratura statistica, alle quali ne sono state aggiunte dagli autori altre tre. Per definire la rappresentatività occorre individuare il contesto di riferimento. Una accezione molto diffusa di rappresentatività del campione è che esso sia una “miniatura” della popolazione stessa. Ovviamente è utopistico pensare di raggiungere questo obiettivo per tutte le variabili, più realisticamente ciò avviene solo per alcune variabili note (ad esempio il genere e l'età degli individui). Naturalmente l'aspettativa è quella di realizzare in tutto o in parte la rappresentatività di un campione anche rispetto alle variabili di indagine. Tale aspettativa dipende però dal legame tra le variabili rispetto alle quali il campione è rappresentativo e quelle oggetto di indagine e non sempre viene realizzata.

Da quanto detto discende che un campione casuale può essere probabilistico e non necessariamente rappresentativo. Viceversa l'essere probabilistico per un campione implica la casualità ma non la rappresentatività, mentre un campione rappresentativo può essere o non essere casuale e/o probabilistico. Il ricercatore sa se il campione è casuale e probabilistico, mentre non sa se è rappresentativo, se non facendo riferimento a esperienze passate o a particolari verifiche a posteriori, ad esempio con l'analisi delle sub-popolazioni, come si avrà occasione di esaminare nel Cap. 3. L'assenza dei primi due attributi implica decisioni sulla scelta dell'approccio inferenziale, come emergerà nel Cap. 3, mentre l'assenza di rappresentatività<sup>2</sup> potrebbe indicare la presenza

---

<sup>2</sup> Negli ultimi anni è stato sviluppato in letteratura un nuovo concetto di rappresentatività legato al contesto di qualità dell'indagine. È opinione diffusa che non è il tasso di non risposta a rendere le stime distorte quanto le differenze tra rispondenti e non nei confronti della variabile di interesse (J. Bethlehem et al.

degli errori non campionari, che devono quindi essere individuati e corretti con le metodologie proposte nei Capp. 4, 5 e 6.

### 1.2.3 Fasi di un'indagine campionaria

L'indagine statistica, sia essa censuaria o campionaria, è una procedura complessa che si articola in una successione di fasi di seguito elencate:

1. *Formalizzazione degli obiettivi, individuazione della popolazione obiettivo, definizione spazio-temporale.* Formalizzare gli obiettivi significa individuare le variabili che verranno osservate nella popolazione di riferimento. Il numero di variabili oggetto della ricerca dipende dalla complessità degli obiettivi: un'indagine sulla lettura dei quotidiani è molto più semplice dell'indagine sul tempo libero, che coinvolge un più ampio numero di variabili rispetto alla prima. La definizione già fornita di popolazione target deve essere circoscritta nel tempo e nello spazio. Ad esempio, nell'indagine sulla lettura dei quotidiani la popolazione obiettivo potrebbe essere costituita dai residenti con più di 15 anni (i residenti con meno di 15 anni potrebbero essere scarsamente interessati alla lettura dei quotidiani) in una definita area territoriale (comune, provincia, ecc.) in un fissato arco temporale (settimana, mese, anno, ecc.).
2. *Costruzione del frame.* Il frame deve essere costruito con cura in modo che contenga tutte le unità della popolazione target. Ad ogni unità del frame corrisponde una etichetta che identifica l'unità di rilevazione, il cui ordinamento può essere casuale o non casuale. Ad esempio, l'ufficio del personale di una azienda può costruire un frame dei suoi dipendenti utilizzando l'ordine alfabetico del cognome (in questo caso l'ordinamento delle unità della lista può ritenersi casuale rispetto ad una usuale variabile d'indagine) ovvero può elencare i dipendenti in funzione della data di assunzione (in questo caso l'ordinamento non è casuale, ma in funzione del tempo, e può essere un vantaggio se la variabile oggetto d'indagine è sensibile a tale ordinamento come, ad esempio, l'avanzamento di carriera o gli scatti di stipendio). La costruzione del frame implica anche l'inserimento, per ogni unità, di tutte le informazioni ausiliarie che possono essere utili per la estrazione del campione o la costruzione dello stimatore; ad esempio, per quanto riguarda l'elenco dei dipendenti tali informazioni potrebbero essere: il genere, il titolo di studio, il livello attribuito al momento dell'assunzione, il livello attuale, lo stipendio, ecc.

---

2008; 7<sup>mo</sup> Programma Quadro dell'UE e i Work Package 1, 2, 3, 4, 5, 6 a cura di N. Shlomo et al. 2009). Pertanto, il campione dei rispondenti potrebbe essere ancora rappresentativo del campione programmato e le stime calcolate su tale campione potrebbero essere ancora realistiche. Per sostenere questa tesi occorre quindi individuare degli indicatori di rappresentatività del campione osservato, che sono stati chiamati R-indicators.

3. *Scelta del piano di campionamento e dello schema di estrazione.* Piano di campionamento e schema di estrazione vengono descritti nel Par. 2.2. Tuttavia, lo schema, proprio perché indica il modo con cui le unità vengono selezionate, si ritiene fondamentale per definire il campione probabilistico o meno. Ad esempio, nel piano di campionamento stratificato la popolazione viene suddivisa, come è noto, in sottoinsiemi (strati) omogenei e disgiunti, per ciascuno viene costruito il frame e individuato uno schema di estrazione, in base al quale in ogni strato viene selezionato un campione. Se il metodo di scelta prevede un esperimento casuale con probabilità di selezione costante per ogni unità in ognuno degli strati, si ha un campione probabilistico e l'unione di questi campioni definisce il *piano di campionamento casuale stratificato* (si veda Par. 2.5.7). Al contrario se non è previsto alcun esperimento casuale per la selezione, che viene demandata all'intervistatore, il campione di strato è non casuale e non probabilistico e l'unione di tali campioni porta al *campionamento per quote* (si veda Par. 1.4).
4. *Rilevazione dei dati.* I metodi di rilevazione o raccolta dei dati sono diversi e vengono scelti in relazione ai differenti contesti in cui l'indagine si svolge. Ad esempio, nelle indagini sugli individui e le famiglie, la raccolta dati richiede spesso l'intervento di un intervistatore. In questo ruolo, il rilevatore ha il compito di contattare il possibile rispondente e di somministrare il questionario. Le modalità di rilevazione possono essere diverse in relazione alle risorse disponibili per l'indagine, ai tempi richiesti per la sua realizzazione e alla natura della popolazione obiettivo da studiare. La rilevazione può essere effettuata con interviste faccia a faccia compilando un questionario cartaceo (modalità PAPI – *Paper and Pencil Interviewing*) oppure un questionario elettronico (modalità CAPI – *Computer Assisted Personal Interviewing*), con interviste telefoniche (modalità CATI – *Computer Assisted Telephone Interviewing*) oppure via Internet, proponendo al soggetto l'auto-compilazione del questionario via Web (modalità CAWI – *Computer Assisted Web Interviewing*). In taluni casi l'auto compilazione può essere proposta anche inviando il questionario cartaceo via posta ordinaria.
5. *Analisi dei dati.* I dati campionari vengono utilizzati in una logica inferenziale per stimare i parametri delle variabili di interesse nella popolazione target. In un'indagine ideale, cioè senza alcun tipo di errore, viene utilizzato prevalentemente l'approccio inferenziale *design-based* (Par. 3.1). Tuttavia l'indagine senza alcun tipo di errore difficilmente si realizza; gli argomenti dei capitoli che seguono sono rivolti proprio ai metodi inferenziali per le indagini in presenza di errori. Comunque, una regola che deve essere sempre rispettata nell'analisi delle variabili oggetto di studio è quella di considerare, congiuntamente alle stime, l'errore quadratico medio, l'errore standard, nonché, quando possibile, l'intervallo di confidenza.
6. *Diffusione dei risultati.* Una volta raccolti i dati ed effettuate le analisi, inizia la fase di pubblicizzazione e diffusione dei risultati, che avviene principalmente tramite la diffusione del cosiddetto rapporto di indagine. Esso contiene la descrizione delle fasi d'indagine e l'illustrazione dei principali

risultati, con attenzione alla rappresentatività del campione e all'attendibilità delle stime. In tale fase, il soggetto responsabile dell'indagine svolge un'attività di promozione, diffusione e trasferimento dei risultati descritti nel rapporto con un programma di iniziative finalizzato a raggiungere tutti gli interessati.

Si fa notare che nell'indagine *censuaria* non è prevista la fase 3; mentre in quella *campionaria* le fasi 5 e 6 sono presenti con una logica diversa rispetto a quella dell'indagine globale.

### 1.2.4 Gli errori di un'indagine campionaria

Nelle fasi dell'indagine possono verificarsi delle imperfezioni, alcune con conseguenze molto gravi, altre meno, che in entrambi i casi richiedono interventi correttivi. In generale le imperfezioni determinano delle divergenze, chiamate comunemente *errori*, tra quello che “si sarebbe teoricamente dovuto osservare” e quello che “realmente è stato osservato”. Diverse sono le tipologie di errore (Lessler e Kalsbeek 1992), tuttavia ciò che mina la qualità della ricerca non è il tipo di errore ma l'entità del medesimo. L'inferenza statistica insegna che gli errori di piccole dimensioni sono accettabili, non lo sono se di elevate dimensioni. La linea di demarcazione dell'accettabilità non può essere individuata in modo oggettivo, in quanto varia in funzione di diversi fattori come l'obiettivo dell'indagine, gli strumenti di rilevazione utilizzati, il livello di precisione richiesto, i tempi e i costi imposti; ed è per questo che tale linea spesso viene suggerita dall'esperienza e dalla sensibilità del ricercatore.

In genere, nelle indagini si è consapevoli della presenza dell'errore e della sua origine e la letteratura propone diversi metodi di correzione. Nel seguito sono brevemente elencate (mantenendo la stessa numerazione delle fasi) le anomalie che si possono verificare nelle fasi dell'indagine con le loro conseguenze, rinviando ai capitoli successivi i metodi di correzione:

1. La popolazione obiettivo non è specificata, obiettivi non chiari, variabili non coerenti con gli obiettivi: sono errori molto gravi con pesanti conseguenze perché potrebbero portare ad invalidare la ricerca. Non esistono metodi generali di correzione.
2. La non corrispondenza tra la popolazione frame e la popolazione target genera gli *errori di copertura*, che saranno trattati nel Cap. 4.
3. La scelta di procedere alla selezione di un campione implica a priori la rinuncia a conoscere il reale valore del parametro  $T$  e ad accettarne una stima  $\hat{T}$ . Conseguenza di tale scelta è la presenza dell'*errore campionario*, definito come la differenza tra la stima e il reale valore. Per ridurre questo errore si può incrementare la dimensione del campione e/o scegliere un piano di campionamento più efficiente di un altro (si veda Par. 2.7). Nell'ambito dell'inferenza *design-based* è possibile conoscere la dimensione dell'errore campionario solo se il campione è probabilistico.

4. In fase di rilevazione sono frequenti gli *errori di non risposta* che consistono nella mancata partecipazione all'indagine o nel non rispondere ad alcune domande del questionario. Nel primo caso la non risposta è *globale* o *totale*, nel secondo è *parziale*; tali errori saranno ampiamente trattati rispettivamente nei Capp. 5 e 6. La non risposta dipende da molti fattori tra cui un senso di diffidenza nei confronti delle indagini, l'obiettivo della ricerca che può coinvolgere variabili sensibili, il metodo di rilevazione e la comprensione del questionario. In questa fase si possono osservare anche gli *errori di misurazione*, (differenze tra il valore reale della variabile  $y$  e il valore rilevato); di questi alcuni sono imputabili all'intervistatore, altri all'intervistato e altri ancora all'imputazione o alla codifica dei dati.
5. È questa la fase più compromessa dalla presenza degli errori di cui sopra. Infatti, l'analisi dei dati di un'indagine campionaria riguarda le stime che vengono calcolate per i parametri descrittivi delle variabili di interesse, per quelli che ne esprimono le relazioni o per l'analisi delle variabili stesse. È evidente che la presenza degli errori provoca allontanamenti incontrollati delle stime dai rispettivi parametri.
6. L'errore che si commette in questa fase consiste nella non esatta descrizione dell'indagine e dei suoi risultati.

Gli errori descritti possono essere classificati in modi diversi; ad esempio:

- Se si vuole evidenziare l'errore della stima, si distingue l'*errore campionario* dagli *errori non campionari*. Il primo è presente solo nelle indagini campionarie e, come già detto, può essere tenuto sotto controllo in vari modi; i secondi sono gli errori di copertura, di non risposta e di misura che sono presenti nelle indagini campionarie ma anche in quelle censuarie; inoltre, in queste ultime la loro presenza è più rilevante in quanto tali errori aumentano con l'aumentare del numero delle osservazioni. Nei capitoli seguenti saranno proposti i metodi per correggere gli errori di copertura e di non risposta (globale e parziale), mentre non saranno considerati i metodi per correggere gli errori di misura, in quanto nella pratica delle indagini campionarie sono scarsamente utilizzati. Tuttavia, poiché questi errori possono distorcere anche sensibilmente le stime, alla fine del Cap. 3 vengono analizzati gli effetti che producono, per i quali si suggerisce la prevenzione come unico approccio valido di contrasto.
- Se si vogliono distinguere gli errori in base al fatto di avere o meno osservato le unità statistiche oggetto dell'indagine, si distinguono gli errori da *non osservazione* da quelli da *osservazione* (Groves 1989). I primi nascono perché non sono state osservate tutte le unità della popolazione, sono tali l'errore campionario, quello di sotto-copertura e di non risposta; i secondi sono tutti gli errori di misurazione che possono essere generati da chi, a vario titolo, si trova coinvolto nell'indagine, come ad esempio gli intervistatori e gli intervistati, o gli errori generati dal metodo di raccolta dei dati o dalla costruzione del database, ecc.