TEORIA E TECNICA DELL'INDAGINE STATISTICA E DEL CAMPIONAMENTO (MATR.DISPARI)

PRINCIPI DI CAMPIONAMENTO

siamo arrivati fino al piano di pubblicazione ci manca il piano di campionamento: cioe insieme di procedure che si adottano quando facciamo un indagine campionaria piuttosto che esaustiva, esaustiva se è convolta tuttal la pop.

es: esami sono domande a campione = indagine campionaria le domande debono essere diversificate cosi rapp la realtà

NON voglio trarre conclusioni dal campione MA dalla pop tramite il campione



Università degli Studi di Padova



manuela.scioni@unipd.it

Dipartimento di Scienze Statistiche

MANUELA SCIONI

1

da adesso in poi popolazione = pop obiettivo e pop statistica quindi suppongo che le 2 pop coincidono da adesso

PROGRAMMA

— non distinguo piu,pop in campionamento parlo di pop statistica

quando tuttel le unit facenti parte della pop. o meglio della mia lista hanno prob non nulla cioe tutti hanno lapossibilita di entrare nella campione

CAMPIONAMENTI PROBABILISTICI

- Campionamento casuale semplice
- Campionamento stratificato «
- Campionamento a stadi

TUTTI 3 CASUALI

ALTRI TIPI DI CAMPIONAMENTO

- Campionamento sistematico (pseudocasuale)
- Campionamento per quote (non probabilistico) scuole allora selz 5 scuole poi mi appproccio alla

unita selez casualmente dalla lista

selez casuale al interno di ogni stratto, sono indipendenti es: laureati padova gli suddivido rispetto alla variabile faccio una selezione casuale al interno di questo gruppo

le unit stat sono aggregate in un unita meno ampie es: laureati apparrtengono universita corsi di lauree che app a univ che app a atenee, seleziono casualmente le unita aggregante es: corsi di laurea,utile quando non abb l intera lista della pop, es: indagine sulle scuole superiori non sapp ogii i nomi cognomi di tutti pero sappiamo tutte le singola scuola dalla lista degli studenti quindi senza avere la lista di partenza sono riuscioto a fare un campionamento

guardare esempio lezione 26 aprile

CAMPIONAMENTO PROBABILISTICO

- Punto di partenza è P, una popolazione di dimensione N
- c = campione di dimensione n tratto dalla popolazione P
- Campionamento probabilistico:

- coef binomiale
- Definire l'insieme dei possibili campioni $C = \{c_1, c_2, ..., c_M\}$
- Assegnare ad ogni possibile campione una probabilitá p(c)
- Ad ogni unità è associata una probabilità $\pi_i > 0$ di far parte del campione

NON FAREMO LA PROCEDURA cmq LA VEDIAMO PRIMA LEZIONE

 $\pi_i = \sum P(c_j | i \in c_j)$

Definire una procedura per ottenere un campione c con probabilità p(c)

in maiuscolo = popolazione minuscolo= campione

t.c la somma = 1 e ogni unita ha prob non nulla di essrre estratta

> cmq puo avere zero ma deve avere la possibilita di entrare

La distribuzione di probabilità p(c) definisce il disegno di campionamento

esempio slide a mano controllare:

fattoriale; 6 numeri di poss campioni che posso

estarare da 4 sensza rienserimento

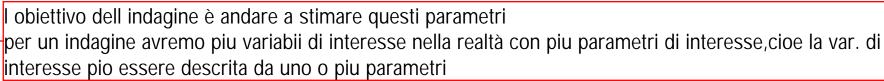
CAMPIONAMENTO DA POPOLAZIONI FINITE

- P è una popolazione di numerosità finita N
- Di conseguenza, il numero di campioni che possiamo costruire è finito, e ad ognuno di essi è associata una probabilità p(c)
- Lo stimatore è funzione dei valori Y; di un campione c:

$$\hat{\theta}_c = f(Y_i, i \in c)$$

parleremo spesso di stimatore della media e totale il camp deve rapp lo pop cosi lo stimatore ha senso, posso associare alcuni parametri allo stimatore che descrivono la distrubuzione

 La distribuzione dello stimatore è composta da un insieme finito di probabilità, ed è perciò assimilabile ad una funzione di probabilità di v.c. discrete



POPOLAZIONE E PARAMETRI (IGNOTI)

P=Popolazione di N elementi

Y=caratteristica di P che vogliamo studiare

Lo studio di Y avviene mediante la stima dei suoi parametri:

- Media
- Totale
- Mediana
- Percentili

 θ è il parametro relativo a Y, nella popolazione P, che intendiamo stimare con l'indagine campionaria

θ Parametro ignoto

$$Media: \bar{Y} = \sum_{\substack{i=1\\N}}^{N} \frac{Y_i}{N}$$

Totale:
$$Y = \sum_{i=1}^{N} Y_i$$

Varianza:
$$S^2 = \frac{1}{N-1} \sum_{\substack{i=1 \ N}}^{N} (Y_i - \bar{Y})^2; \ \sigma^2 = \frac{1}{N} \sum_{i=1}^{N} (Y_i - \bar{Y})^2$$

Covarianza:
$$\sigma_{XY} = \frac{1}{N} \sum_{i=1}^{N} (X_i - \bar{X})(Y_i - \bar{Y})$$

Coefficiente di correlazione:
$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

Se Y è dicotomica (Y = 0 assenza, Y = 1 presenza):
$$\bar{Y} = \text{proporzione} = P = \frac{N_1}{N} \sigma^2 = P \cdot (1 - P)$$

in realta non faremo mai cosi, ma prima estraiamo un campione e poi vediamo errore

PARAMETRI DELLO STIMATORE

esempio altezza media lezione 26 aprile

Distribuzione campionaria:	$P(\hat{\vartheta} = k) = \sum_{c:\widehat{\theta}_c = k} p(c)$
Valore atteso:	$\nabla \hat{a} = \nabla \hat{a} = 0$
	$E(\hat{\theta}) = \sum_{c} \hat{\theta}_{c} p(c)$ ci serve per capire se distorto
Distorsione:	$B(\hat{\theta}) = E(\hat{\theta}) - \theta$
	$D(\theta) - E(\theta) - \theta$
Varianza:	$V_{cm}(\hat{0}) = \sum_{i} (\hat{0} - F_i(\hat{0}))^2 m(s)$
	$Var(\hat{\theta}) = \sum_{c} (\hat{\theta}_{c} - E(\hat{\theta}))^{2} p(c)$
Errore Quadratico Medio:	$MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2 =$
	$\sum_{c} (\hat{\theta}_{c} - \theta)^{2} p(c) = Var(\hat{\theta}) + (B(\hat{\theta}))^{2}$
	C

entrambi gli errori danno varianza di stima piu varianza di rilevazione, quello calc nel esempio era varianza di stima, in nostro errore camp deve essere il piu piccolo possibile è questo I obiettivo di tutte le procedure possibili, errore camp err non camp vanno controllati sempre, di solito quello di rilevazione è piu alto di solito, es: xke le persone non dichiareano il vero es: quando berlusconi vinsse

COMPONENTI DELL'ERRORE

vanno contrallati tutti e due gli errori

xke indagine camp

ERRORE CAMPIONARIO

- Al cambiare del campione cambia la stima di θ , quindi varia $\hat{\theta}$:
- Varianza di stima:

$$Var_S(\hat{\theta}) = \sum_{c}^{N^*} \{\theta_c - E(\hat{\theta})\}^2 p_c$$

Errore campionario di stima

$$Var_S(\hat{\theta})$$

Errore quadratico medio di stima:

$$Mse_S(\hat{\theta}) = \sum_c (\hat{\theta}_c - \theta)^2 p(c) = Var(\hat{\theta}) + (B(\hat{\theta}))^2$$

errorei per errore del rispondente o nostro,, puo anche essere sistematico

ERRORE DI RILEVAZIONE

- È dovuto a errate dichiarazioni (volontarie, dimenticanze, interpretazioni non corrette da parte dell'intervistatore, ...)
- $f \cdot$ Anch'esso contribuisce alla varian ${f z}$ a di $\hat{m heta}$

varianza di rilevazione, in un infagine camp non posso calcolarala

$$Var(\hat{\theta}) = Var_S(\hat{\theta}) + Var_R(\hat{\theta})$$

varianza stima calcolata prima nel esempio 40 min