

Bortfaldsproblemet i stikprøveundersøgelser

Af Carsten Stig Poulsen*)

Indledning

I en artikel i dette tidsskrift¹⁾ rejser Flemming Hansen spørgsmålet om den eksisterende statistiske teoris praktiske anvendelighed på stikprøveundersøgelser af populationer, der består af personer, grupper af personer m.v. Udover den usikkerhed, som følger af, at det kun er en del af populationen, der undersøges, optræder der her tillige kilder til fejl i form af målefejl og bortfald (non-response), som stort set ikke indgår i de teoretisk-statistiske modeller.

Det er hensigten i det følgende at opstille en formel model, som kan belyse nogle sider af bortfaldsproblemet betydning for statistisk inferens. De spørgsmål, som primært behandles er: Hvilken indflydelse har bortfaldet på egenskaberne hos de estimatorer, som traditionelt anvendes i stikprøveundersøgelser, og hvorledes kan undersøgelserne tilrettelægges under hensyntagen til bortfaldet.

Fremstillingen prætenderer ikke at være fuldt dækkende, men repræsenterer nogle indledende overvejelser, som måske kan danne grundlag for videre forskning inden for området.

*) Cand. polit., adjunkt, AUC.

¹⁾ Nr. 4/1977.

Bortfaldsproblemet i stikprøveundersøgelser

Af Carsten Stig Poulsen*)

Indledning

I en artikel i dette tidsskrift¹⁾ rejser Flemming Hansen spørgsmålet om den eksisterende statistiske teoris praktiske anvendelighed på stikprøveundersøgelser af populationer, der består af personer, grupper af personer m.v. Udover den usikkerhed, som følger af, at det kun er en del af populationen, der undersøges, optræder der her tillige kilder til fejl i form af målefejl og bortfald (non-response), som stort set ikke indgår i de teoretisk-statistiske modeller.

Det er hensigten i det følgende at opstille en formel model, som kan belyse nogle sider af bortfaldsproblemet betydning for statistisk inferens. De spørgsmål, som primært behandles er: Hvilken indflydelse har bortfaldet på egenskaberne hos de estimatorer, som traditionelt anvendes i stikprøveundersøgelser, og hvorledes kan undersøgelserne tilrettelægges under hensyntagen til bortfaldet.

Fremstillingen prætenderer ikke at være fuldt dækkende, men repræsenterer nogle indledende overvejelser, som måske kan danne grundlag for videre forskning inden for området.

*) Cand. polit., adjunkt, AUC.

¹⁾ Nr. 4/1977.

Bortfaldet som problem

Ved *bortfald* vil vi i det følgende forstå manglende information om nogle af de analyseenheder, som var planlagt at skulle indgå i undersøgelsen. Denne fejlkilde knytter sig ikke specielt til stikprøveundersøgelser, men kan lige såvel optræde ved totaltællinger.

Bortfaldet som problem forstærkes i det omfang de bortfaldne analyseenheder adskiller sig væsentligt fra resten af stikprøven med hensyn til det undersøgte kendetegn. Da vil de skøn over forholdene i populationen, som bygger på den gennemførte stikprøve, nødvendigvis blive fejlagtige. I forbindelse med interviewundersøgelser kan der sondres mellem to hovedtyper af bortfald: respondenter, som ikke træffes og nægttere. De to typer af bortfald adskiller sig på en række punkter. Man vil således ved genbesøg være i stand til at nedbringe antallet af ikke-trufne respondenter, mens det er mere krævende at overtale folk med nægter-indstilling til at deltage i undersøgelsen. Forudsætningen om, at bortfaldet ikke afviger væsentligt fra resten af populationen kan også være mere hasaderet i forbindelse med nægtergruppen end for gruppen af ikke-trufne. Endelig skal man være opmærksom på en mulig sammenhæng mellem formålet med undersøgelsen og tilbøjeligheden til bortfald.

Hvis man f.eks. udtager en stikprøve med henblik på en vurdering af antallet af sortseere, er det ikke utænkeligt, at nægterprocenten blandt sortseerne er større end blandt licensbetalere. (En anden mulig fejlkilde er her målefejl i form af urigtige oplysninger).

Vi skal overalt i det følgende antage, at bortfaldet er uden sammenhæng med det undersøgte kendetegn.

En statistisk model med bortfald

Problemstillingen i en model med bortfald sammenlignet med traditionel statistisk stikprøvetæori er søgt anskueliggjort i figur 1.

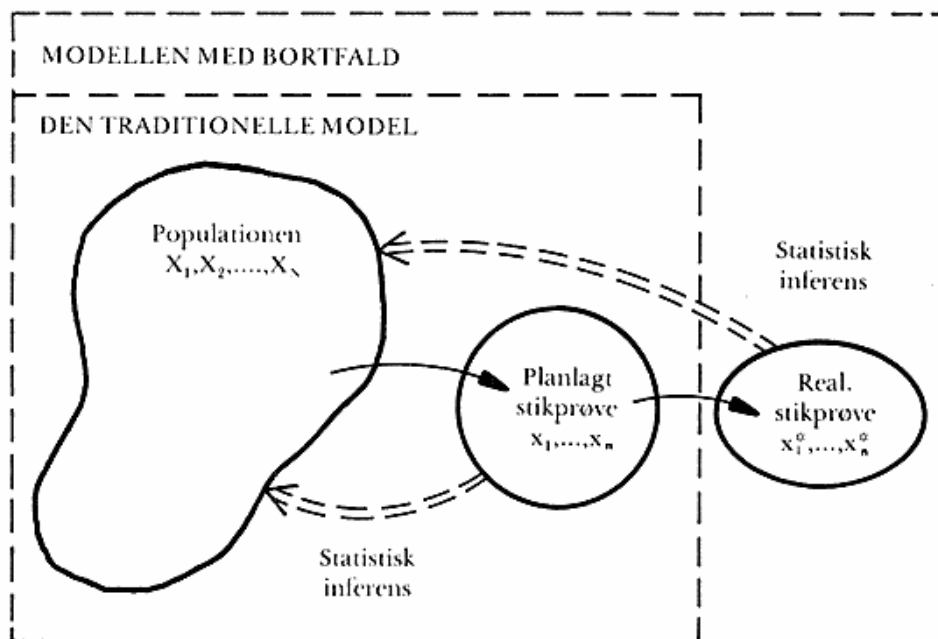


Fig. 1.

Grundlaget i den traditionelle model er en stikprøve, som er udtaget *tilfældigt*, d.v.s. sandsynligheden for, at en vilkårlig enhed i populationen indgår i stikprøven er *kendt* – men ikke nødvendigvis den samme for alle enheder.

Derved etableres en sandsynlighedsteoretisk sammenhæng mellem kendetegn i populationen og observerede kendetegn i stikprøven, som danner grundlag for *statistisk inferens* med angivelse af den dertil knyttede usikkerhed. Udarbejdelsen af en stikprøveplan bygger på et *effektivitetskriterium*: For et givet økonomisk budget søges stikprøveusikkerheden minimeret eller – omvendt – en given usikkerhed ønskes opnået ved de lavest mulige omkostninger.

I den mere generelle model med bortfald reduceres den planlagte stikprøve, som følge af non-response hos nogle af de udvalgte enheder. Stikprøven omfatter kun n^* enheder i stedet for de planlagte n , et bortfald på $n - n^*$. Problemet er nu, i hvilket omfang inferens reglerne fra den traditionelle model kan anvendes i dette mere generelle tilfælde.

I den model, der opstilles for at belyse dette spørgsmål, anvendes følgende notation:

X_v : Værdien af kendetegnet for v 'te enhed i populationen,
 $v = 1, 2, \dots, N$.

x_i : Værdien af kendetegnet for i 'te enhed i den planlagte stikprøve,
 $i = 1, 2, \dots, n$.

x_j^* : Værdien af kendetegnet for j 'te enhed i den realiserede stikprøve, $j = 1, 2, \dots, n^*$.

θ_v : Sandsynligheden for, at v 'te enhed i populationen udtages til den planlagte stikprøve, $v = 1, 2, \dots, N$. $0 < \theta_v < 1$.

ρ_v : Sandsynligheden for, at v 'te enhed i populationen *ikke* indgår i den realiserede stikprøve, skønt den er udtaget til den planlagte stikprøve, $v = 1, 2, \dots, N$. $0 < \rho_v < 1$.

Vi antager, at formålet med stikprøven er at estimere middeltallet i populationen:

$$\bar{X} = \frac{1}{N} \sum_{v=1}^N X_v \quad (1)$$

Simpel tilfældig udvælgelse

Det forudsættes, at den planlagte stikprøve udtages ved *simpel tilfældig udvælgelse*. Dette indebærer, at alle enheder i populationen har samme sandsynlighed for at blive valgt:

$$\theta_v = \frac{n}{N} ; v = 1, 2, \dots, N \quad (2)$$

Da gælder følgende sætning i *den traditionelle model*:

Sætning 1: Middeltallet i stikprøven

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (3)$$

er en *middelret* estimator over middeltallet i populationen \bar{X} .

Bevis: Lad indikatorvariablen α_{ν} være defineret ved

$$\alpha_{\nu} = \begin{cases} 1, & \text{hvis } X_{\nu} \text{ indgår i stikprøven} \\ 0, & \text{ellers} \end{cases} \quad (4)$$

Da kan estimatoren \bar{x} skrives

$$\bar{x} = \frac{1}{n} \sum_{\nu=1}^N \alpha_{\nu} \cdot X_{\nu} \quad (5)$$

Den forventede værdi af \bar{x} , $E(\bar{x})$, kan findes:

$$E(\bar{x}) = \frac{1}{n} \sum_{\nu=1}^N E(\alpha_{\nu}) \cdot X_{\nu}. \quad (6)$$

Idet

$$E(\alpha_{\nu}) = 1 \cdot \Pr\{\alpha_{\nu}=1\} + 0 \cdot \Pr\{\alpha_{\nu} = 0\} = \frac{n}{N}$$

fås:

$$E(\bar{x}) = \frac{1}{n} \sum \frac{n}{N} X_{\nu} = \bar{X} \quad \blacktriangleleft$$

Uden bevis²⁾ anføres.

Sætning 2: Den teoretiske varians på estimatoren \bar{x} er givet ved

$$V(\bar{x}) = \frac{\sigma^2}{n} \frac{N-n}{N-1} \approx \sigma^2 \left(\frac{1}{n} - \frac{1}{N} \right) \quad (7)$$

hvor

$$\sigma^2 = \frac{1}{N} \sum_{\nu=1}^N (X_{\nu} - \bar{X})^2 \quad (8)$$

er variansen i populationen.

²⁾ Et bevis kan f.eks. findes hos Jensen (1974).

For $N \rightarrow \infty$ fremkommer den velkendte formel for variansen på middeltallet af uafhængige observationer fra en uendelig population.

I modellen med bortfald skal vi nu betragte

$$\bar{x}^* = \frac{1}{n^*} \sum_{j=1}^{n^*} x_j^* \quad (9)$$

d.v.s. middeltallet i den gennemførte stikprøve som estimator for \bar{X} . Med α_v defineret ved (4) kan (9) skrives:

$$\bar{x}^* = \frac{1}{n^*} \sum_{v=1}^N \alpha_v X_v \quad (10)$$

Her er fordelingen af α_v imidlertid ændret. Vi antager, at sandsynligheden for at indgå i den realiserede stikprøve kan skrives:

$$\Pr\{\alpha_v = 1\} = \theta_v \cdot (1 - \rho_v) = \frac{n}{N} (1 - \rho_v) \quad (11)$$

d.v.s. som produktet af sandsynligheden for at blive udtaget til den planlagte stikprøve og tilbøjeligheden (sandsynligheden) for at indgå i den realiserede stikprøve. Man forestiller sig altså, at populationen ikke blot er karakteriseret ved kendetegnet X_v , men også ved bortfaldstilbøjelighederne ρ_v . En enhed består af et talpar (X_v, ρ_v) , som ved udtrækning repræsenteres ved (x_i, ρ_i) . x_i vil da indgå i den realiserede stikprøve med sandsynligheden $1 - \rho_i$.

Mens θ_v alene afhænger af den valgte stikprøveplan, er ρ_v i vidt omfang uden for planlæggerens kontrol.

For $\rho_v = 0$ fremkommer den traditionelle model som et specialtilfælde.

Vi skal nu vise følgende:

*Sætning 1**: Middeltallet \bar{x}^* i den realiserede stikprøve er som hovedregel ikke en middeltal estimator for middeltallet i \bar{X} populationen.

Bevis: Med α_v defineret ved (4) kan antallet af enheder n^* i den gennemførte stikprøve skrives:

$$n^* = \sum_{v=1}^N \alpha_v$$

d.v.s. n^* er en stokastisk variabel. Den forventede værdi af \bar{x}^* er givet som:

$$E(\bar{x}^*) = E\left(\frac{1}{n^*} \sum_{v=1}^N \alpha_v \cdot X_v\right) = \sum_{v=1}^N E\left(\frac{\alpha_v}{n^*}\right) \cdot X_v \quad (12)$$

\bar{x}^* vil derfor kun være middelværdi, såfremt

$$\sum_{v=1}^N E\left(\frac{\alpha_v}{n^*}\right) \cdot X_v = \frac{1}{N} \sum_{v=1}^N X_v$$

$$\Leftrightarrow \sum_{v=1}^N \left(E\left(\frac{\alpha_v}{n^*}\right) - \frac{1}{N}\right) \cdot X_v = 0 \quad (13)$$

hvilket ikke kan antages at gælde i almindelighed. ◀

Vi skal nu undersøge (13) i nogle specialtilfælde. Hvis n^* er rimelig stor, vil α_v kunne betragtes som uafhængig af $n^* = \sum \alpha_v$. Da gælder:

$$E\left(\frac{\alpha_v}{n^*}\right) = E\left(\frac{1}{n^*}\right) \cdot E(\alpha_v)$$

Anvendelse af (11) giver

$$E(\alpha_v) = \frac{n}{N} (1 - \rho_v)$$

og dermed kan betingelsen (13) skrives

$$E\left(\frac{1}{n^*}\right) \cdot \frac{n}{N} \sum_{v=1}^N (1 - \rho_v) \cdot X_v - \bar{X} = 0 \quad (14)$$

Introduceres nu mere formelt antagelsen om *uafhængighed mellem bortfaldssandsynlighed og værdien af kendetegnet*, haves

$$\sum_{v=1}^N (\rho_v - \bar{\rho}) (X_v - \bar{X}) = 0$$

$$\Rightarrow \sum_{v=1}^N (1 - \rho_v) X_v = N \cdot \bar{X} (1 - \bar{\rho})$$

som indsat i (14) giver

$$E\left(\frac{1}{n^*}\right) = \frac{1}{n(1-\bar{\rho})} \quad (15)$$

Da $1/n^*$ er en konveks funktion af n^* gælder det ifølge Jensens ulighed, at

$$E\left(\frac{1}{n^*}\right) > \frac{1}{E(n^*)}$$

Vi har tillige

$$E(n^*) = \sum_{v=1}^N E(\alpha_v) = \sum_{v=1}^N \frac{n}{N} (1 - \rho_v) = n(1 - \bar{\rho})$$

og dermed

$$E\left(\frac{1}{n^*}\right) > \frac{1}{n(1-\bar{\rho})}$$

Selv når n^* er stor gælder sætning 1* stadig. Variationen i det gennemførte antal observationer er medvirkende til, at \bar{x}^* er et skævt skøn over \bar{X} . Skævheden vokser med variansen på n^* .

Vi ser dernæst på *en vigtig undtagelse* fra hovedreglen i sætning 1*:

*Korollar 1**: Såfremt bortfaldssandsynligheden ρ_v er konstant, $\rho_v = \bar{\rho}$, vil \bar{x}^* være en middelret estimator for \bar{X} .

Bevis: Når $\rho_v = \bar{\rho}$, er $\Pr\{\alpha_v = 1\} = \frac{n}{N}(1 - \bar{\rho})$ konstant, uafhængig af v . Antages der tillige uafhængighed mellem enhederne, hvad angår deres tilstedeværelse i den realiserede stikprøve, vil fordelingen af n^* være

givet ved en binomialfordeling med antalsparameteren N og sandsynlighedsparameteren $\frac{n^*}{N}(1-\bar{p})$. Der gælder nu generelt³⁾:

$$E\left(\frac{\alpha_v}{n^*}\right) = E\left(\frac{1}{n^*} \cdot E(\alpha_v | n^*)\right) \quad (16)$$

hvor $E(\alpha_v | n^*)$ betegner den betingede middelværdi af α_v , givet summen $n^* = \sum \alpha_v$. Da n^* er binomialfordelt gælder⁴⁾:

$$E(\alpha_v | n^*) = \Pr\{\alpha_v=1 | n^*\} = \frac{n^*}{N} \quad (17)$$

$$\Rightarrow E\left(\frac{\alpha_v}{n^*}\right) = E\left(\frac{1}{n^*} \cdot \frac{n^*}{N}\right) = \frac{1}{N}$$

og dermed

$$E(\bar{x}^*) = \frac{1}{N} \sum X_v = \bar{X} \quad \blacktriangleleft$$

En intuitiv forklaring på dette resultat er:

Når bortfaldssandsynligheden er konstant over enhederne i populationen er den »mekanisme«, som udvælger enheder fra den planlagte til den realiserede stikprøve helt analog med udvalgsproceduren fra populationen til den planlagte stikprøve.

Usikkerheden på skønnet \bar{x}^* skal nu vurderes. Som det kunne forventes gælder:

*Sætning 2**: Bortfald bevirker en forøgelse af usikkerheden, målt ved variansen på estimatoren \bar{x}^* .

Delvist bevis: Variansen på \bar{x}^* , $V(\bar{x}^*)$, kan skrives⁵⁾

$$V(\bar{x}^*) = E[V(\bar{x}^* | n^*)] + V[E(\bar{x}^* | n^*)] \quad (18)$$

For givet n^* vil $E(\bar{x}^* | n^*)$ i det generelle tilfælde med ρ_v forskellig afvige fra \bar{X} og afhænge af n^* . Derfor er andet led i (18) positivt. Ligeledes vil den betingede varians $V(\bar{x}^* | n^*)$ være større end variansen i den traditionelle model med samme stikprøvestørrelse:

³⁾ For produktet af to stokastiske variable X og Y gælder det, at $E(XY) = E(XE(Y|X)) = E(YE(X|Y))$, se f.eks. Jensen (1974), p. 55.

⁴⁾ Den betingede sandsynlighed for, at et vilkårligt binomiallovs forsøg i en forsøgsrække på N er »gunstigt«, givet n^* »gunstige« ialt, er n^*/N , uafhængigt af den marginale sandsynlighed for »gunstigt« udfald i et forsøg.

⁵⁾ se f.eks. Jensen (1974), p. 59.

$$V(\bar{X}^* | n^*) > \sigma^2 \left(\frac{1}{n^*} - \frac{1}{N} \right)$$

$$\rightarrow E[V(\bar{X}^* | n^*)] > \sigma^2 \left(E\left(\frac{1}{n^*}\right) - \frac{1}{N} \right)$$

Da $E\left(\frac{1}{n^*}\right) > 1/E(n^*) = 1/n(1-\bar{\rho})$ fås

$$E[V(\bar{X}^* | n^*)] > \sigma^2 \left(\frac{1}{n(1-\bar{\rho})} - \frac{1}{N} \right) > V(\bar{X}) \quad \blacktriangleleft$$

Specialtilfældet med $\rho_V = \bar{\rho}$ giver følgende

*Korollar 2**: Såfremt bortfaldssandsynligheden ρ_V er konstant, $\rho_V = \bar{\rho}$, vil

$$V(\bar{X}^*) \approx \sigma^2 \left(E\left(\frac{1}{n^*}\right) - \frac{1}{N} \right) \quad (19)$$

Intuitivt bevis: Som nævnt er situationen med konstant bortfaldssandsynlighed analog til simpel tilfældig udvælgelse i den traditionelle model. For givet n^* har vi derfor ifølge sætning 1 og 2:

$$E(\bar{X}^* | n^*) = \bar{X}$$

$$V(\bar{X}^* | n^*) = \frac{\sigma^2}{n^*} \left(\frac{N-n^*}{N-1} \right) \approx \sigma^2 \left(\frac{1}{n^*} - \frac{1}{N} \right)$$

og dermed

$$V(\bar{X}^*) \approx \sigma^2 \left(E\left(\frac{1}{n^*}\right) - \frac{1}{N} \right) \quad \blacktriangleleft \quad (20)$$

Udnytter vi, at⁶⁾

$$E\left(\frac{1}{n^*}\right) \approx \frac{1}{E(n^*)} + \frac{V(n^*)}{[E(n^*)]^3} = \frac{1+\gamma_{n^*}^2}{E(n^*)} \quad (21)$$

hvor $\gamma_{n^*}^2 = V(n^*)/[E(n^*)]^2$ er *den relative varians*, kan (20) skrives

⁶⁾ (21) forekommer ved rækkeudvikling af $f(x) = \frac{1}{x}$ omkring $x = E(x)$ og efterfølgende middelværdidannelse.

$$V(\bar{x}^*) \approx \sigma^2 \left(\frac{1+\gamma^2}{n(1-\rho)} \cdot \frac{n^*}{N} - \frac{1}{N} \right) \quad (22)$$

hvoraf det klart fremgår, at usikkerheden i modellen med bortfald øges, ikke blot som følge af det reducerede antal observationer, $n > n(1-\rho)$, men også på grund af variationen i n^* , målt ved den relative varians $\gamma_{n^*}^2$.

Stratificeret tilfældig udvælgelse

Hovedkonklusionen af analysen over bortfaldets virkninger ved simpel tilfældig udvælgelse er, at anvendeligheden af den traditionelle model afhænger af, hvor *homogen* populationen er m.h.t. bortfaldssandsynligheden. Hvor denne forudsætning ikke er opfyldt, fører dette naturligt tanken hen på en stratifikation af populationen med henblik på dannelse af delpopulationer med (tilnærmelsesvis) konstant bortfaldssandsynlighed, hvorfra der udtages en stikprøve ved simpel tilfældig udvælgelse.

Lad populationen være opdelt i L strata, som hver karakteriseres ved

- 1) stratumstørrelsen N_h
- 2) stratummiddeltallet \bar{x}_h
- 3) stratumvariansen σ_h^2
- 4) bortfaldssandsynligheden ρ_h , som er konstant inden for et stratum,

hvor $h = 1, 2, \dots, L$. Fra hvert stratum udvælges ved simpel tilfældig udvælgelse n_h enheder, der efter bortfald reduceres til n_h^* . Da vil $x_h^* \approx X_h$, iflg. det foregående og

$$\bar{x}_{str}^* = \sum_{h=1}^L \frac{N_h}{N} \cdot \bar{x}_h^* ; \quad N = \sum_{h=1}^L N_h \quad (23)$$

vil dermed være en middelfret estimator over populationsgennemsnittet $\bar{X} = \sum (N_h/N) \cdot \bar{X}_h$. Variansen på \bar{x}_{str}^* er givet ved

$$\begin{aligned}
V(\bar{x}_{str}^*) &= \sum_{h=1}^L \frac{N_h^2}{N^2} V(\bar{x}_h^*) = \sum_{h=1}^L \frac{N_h^2}{N^2} \sigma_h^2 \left(E\left(\frac{1}{n_h^*}\right) - \frac{1}{N} \right) \\
&\approx \sum_{h=1}^L \frac{N_h^2}{N^2} \cdot \sigma_h^2 \left(\frac{1+\gamma^2}{n_h(1-\bar{\rho}_h)} - \frac{1}{N} \right) \quad (24)
\end{aligned}$$

jf. (22). Af (24) ses, at stratifikationen bør ske således at strata er *homogene*, hvad angår det kendetegn, der ønskes undersøgt. Da homogenitet med hensyn til bortfald inden for et stratum er forudsætningen for modellen, kan der opstå problemer, hvis disse to kriterier er *konkurrende*. Da må stratifikationen ske ved en afvejning mellem homogenitet m.h.t. kendetegnet og bevarelse af den simple statistiske model. Det er åbenbart, at en sådan afvejning vil være uhyre vanskelig at gennemføre i praksis.

Som det også fremgår af (24) afhænger usikkerheden på skønnet tillige af *stikprovens fordeling på strata*. Dette problem skal nu drøftes nøjere.

Fordelingen på strata af en given stikprøve n i den traditionelle model siges at være *optimal*, hvis variansen $V(\bar{x}_{str})$ minimeres under bibetingelsen $\sum n_h = n$. Det kan vises⁷⁾, at løsningen resulterer i fordelingen

$$n_h = \frac{N_h \sigma_h}{\sum_{h=1}^L N_h \sigma_h} \cdot n \quad (25)$$

I modellen med bortfald kan det tilsvarende problem formuleres med udgangspunkt i (24) og løsningen bestemmes efter samme principper som i den traditionelle model. Antager vi for simpelhedens skyld, at den relative varians er konstant over strata, $\gamma_{n_h^*}^2 = \gamma^2$, kan det vises, at

$$n_h = \frac{N_h \sigma_h / \sqrt{1-\bar{\rho}_h}}{\sum N_h \sigma_h / \sqrt{1-\bar{\rho}_h}} \cdot n \quad (26)$$

⁷⁾ Jensen (1974), p. 124

giver den optimale fordeling af den *planlagte* stikprøve på strata i modellen med bortfald. Det ses, at jo større bortfaldssandsynligheden p_h er, desto større andel af stikprøven skal udtages fra stratum h , men da p_h optræder under kvadratrodstegn, skal forskellene mellem \bar{p}_h være store, før det giver væsentlige udslag i stikprøvens fordeling⁸⁾.

Korrektion for bortfald

I den nævnte artikel af Flemming Hansen diskuteres forskellige metoder til *opvejning* med det primære formål at korrigere for bortfaldet. Den sammenhæng, som herved tilsyneladende etableres mellem bortfaldsproblemet og behovet for opvejning, er uheldig. Vejning er påkrævet, når den udtagne stikprøve ikke er repræsentativ for populationen med hensyn til karakteristika, som har eller antages at have indflydelse på det kendetegn, man ønsker at undersøge. Stikprøvens information om kendetegnet, er *betinget* af disse karakteristika, og for at kunne skønne over kendetegnets udbredelse i *hele* populationen må man veje med vægte, som er bestemt af *de relative stratumstørrelser*, N_h/N . Dette gælder uafhængigt af den måde, hvorpå stikprøven er udvalgt og uanset omfanget af bortfald. Med kendskab til vægtene N_h/N repræsenterer denne vejning ikke noget *statistisk* problem. Er vægtene ikke kendte, kan der være tale om at stratificere den udtagne stikprøve, hvis $n_h^*/n^* \approx N_h/N$, og da indføres endnu en kilde til usikkerhed i den statistiske inferens.

Bortfaldet som statistisk problem bør i princippet inddrages allerede ved tilrettelæggelsen af undersøgelsen. Man kan søge at udnytte information fra tidligere undersøgelser af samme art. Det er ikke utænkeligt, at bortfaldssandsynligheden afhænger af en række generelle kriterier så som alder, køn, erhverv, uddannelse m.m., som hyppigt anvendes ved stratifikation af populationen. Forskelle i bortfaldssandsynligheder *mellem* disse strata kan testes med anvendelse af oplysninger om det faktiske bortfald, $n_h - n_h^*$, i tidligere undersøgelser, som da kan indgå i den planlagte stikprøves fordeling på strata. Man antager her implicit, at p_h er konstant over individer *inden for* hvert stratum.

⁸⁾ Det bemærkes, at $(1-p_h)$ indgår i (26) som grænseomkostningerne i en traditionel model, der optimeres under hensyntagen til en lineær omkostningsfunktion $C = c_0 + \sum_{h=1}^H c_h n_h$, se Jensen (1974), pp. 130-131.

For så vidt angår den del af bortfaldet, der skyldes ikke-trufne respondenter, er det muligt at tage hensyn til bortfaldet ved at forlade den simple tilfældige udvælgelse med $\theta_{\nu} = \frac{n}{N}$ og tillade θ_{ν} at variere. Med kendskab til ρ_{ν} kan θ_{ν} fastsættes således, at sandsynligheden for at indgå i den gennemførte stikprøve, $\Pr\{\alpha_{\nu} = 1\} = \theta_{\nu}(1 - \rho_{\nu})$, er konstant. Herved opnås at den traditionelle model kan anvendes på trods af den tilfældige variation i bortfaldet. Personparameteren ρ_{ν} må også her søges forklaret og estimeret ved en række generelle egenskaber ved enhederne i populationen.

Litteratur:

Hansen, Flemming: Sampling teori og anvendt statistik. Erhvervsøkonomisk Tidsskrift nr. 4, 1977.

Jensen, Niels Erik (red.): Stikprøvetæori. København 1974.

For så vidt angår den del af bortfaldet, der skyldes ikke-trufne respondenter, er det muligt at tage hensyn til bortfaldet ved at forlade den simple tilfældige udvælgelse med $\theta_{\nu} = \frac{n}{N}$ og tillade θ_{ν} at variere. Med kendskab til ρ_{ν} kan θ_{ν} fastsættes således, at sandsynligheden for at indgå i den gennemførte stikprøve, $\Pr\{\alpha_{\nu} = 1\} = \theta_{\nu}(1 - \rho_{\nu})$, er konstant. Herved opnås at den traditionelle model kan anvendes på trods af den tilfældige variation i bortfaldet. Personparameteren ρ_{ν} må også her søges forklaret og estimeret ved en række generelle egenskaber ved enhederne i populationen.

Litteratur:

Hansen, Flemming: Sampling teori og anvendt statistik. Erhvervsøkonomisk Tidsskrift nr. 4, 1977.

Jensen, Niels Erik (red.): Stikprøvetæori. København 1974.