

Kalibreringsrapport

1 Inledning

I en urvalsundersökning är alltid skattningarna behäftade med *urvalsfel* beroende på att endast en delmängd (urval) av populationen studeras. Ett annat fel uppkommer om vi inte lyckas få svar från alla personer (bortfall) och om de avviker från de svarande med avseende på undersökningsvariablerna. Detta fel kallas för *bortfallsfel*.

För att underlätta användningen av statistiken är det värdefullt om storleken på felen kan uppskattas. Av nämnda feltyper är det endast storleken på urvalsfelet som kan skattas med hjälp av urvalsinformation. Kunskap om bortfallsfelet kan i regel bara fås på ett indirekt och approximativt sätt genom att utnyttja registervariabler.

Både urvalsfel och bortfallsfel kan reduceras genom att använda ett effektivt uppräkningsförfarande. I följande avsnitt redovisas hur det görs i denna undersökning.

2 Parametrar

De parametrar som skattas i denna undersökning är totaler och kvoter.

3 Hjälpinformation

Viss hjälpinformation utnyttjas vanligtvis även före estimationen, t.ex. för bildande av stratifierade urvalsdesigner. I denna undersökning drogs ett klusterurval. I första steget drogs ett stratifierat urval skolor. Elever i utvalda skolor totalundersöktes. En enkät skickades därefter till föräldrarna till utvalda elever.

På grund av bortfall i enkätundersökningen så används hjälpinformation på individnivå i avseende att reducera de skevheter som detta bortfall kan leda till.

Det centrala arbetet för att få god kvalitet på skattningarna, då kalibrerings-estimatoren används, är att använda ”stark” hjälpinformation. I följande avsnitt beskrivs detta arbete för denna undersökning.

3.1 Tänkbara hjälpvariabler

Vid val av hjälpvariabler är det tre kriterier som ska beaktas (se Lundström och Särndal 2001):

- (i) Det första kriteriet är att variabeln samvarierar väl med svarsbenägenheten (-sannolikheten). Det är det viktigaste kriteriet eftersom det leder till en minskning av bortfallsskevheter för alla skattningar.

- (ii) Det andra kriteriet är att variabeln samvarierar väl med (viktiga) målvariabler. Om så är fallet minskar bortfallsbiasen för de skattningar som byggs upp av dessa målvariabler. Även variansen minskar för dessa skattningar.
- (iii) Det tredje kriteriet är att variabeln avgränsar (viktiga) redovisningsgrupper. Det leder framförallt till minskad varians i skattningar för dessa redovisningsgrupper.

I den här undersökningen innehåller enkäten frågor av mycket skiftande karaktär. Därför är det främst punkterna (i) och (iii) som kan beaktas, vilket innebär att kalibreringen främst tjänar till att reducera den skevhet som bortfallet troligtvis ger upphov till.

Tänkbara hjälpvariabler, det vill säga variabler som tros uppfylla de ovan uppsatta kriterierna, hämtades ifrån Registret över totalbefolkningen (RTB), Skolregistret och Utbildningsregistret (UREG). En genomgång av variablerna resulterade i att sex variabler valdes ut. De sammanslagningar av kategorier som gjorts baseras på kunskaper från tidigare kalibreringar.

Hjälpvariablerna är definierade enligt tabell 1.

Tabell 1. Tänkbara hjälpvariabler

Variabel (benämning)	Kategorier (koder)
MODERNS FÖDELSELAND	1 = Född i Sverige 2 = Född i övriga Världen
MODERNS ÅLDER	1 = 28-39 2 = 40-45 3 = 46 -
REGION ¹	1 = Storstadskommuner 2 = Övriga kommuner
HÖGSTA UTBILDNING FÖRÄLDRAR	1 = Grundskoleutbildning eller lägre 2 = Gymnasial utbildning 3 = Eftergymnasial utbildning
ELEVENS KÖN	1 = Pojke 2 = Flicka
HUVUDMAN	1 = Kommunal skola 2 = Fristående skola

Med högsta utbildning föräldrar avses högsta avslutade utbildning till och med vårterminen 2010. Utbildningen avser den förälder som har högst utbildning.

Med region avses den region där eleven var folkbokförd i februari 2011. Storstadskommuner är Stockholm, Göteborg och Malmö kommun.

I följande avsnitt analyserar vi variablerna i tabell 1 för att slutligen bestämma en hjälpvektor.

4 Analys av hjälpinformation

4.1 Kriterium 1: Variabeln samvarierar med svarsbenägenheten

För att se huruvida hjälpvariablerna uppfyller det första kriteriet, studeras sambandet mellan den dikotoma variabeln svarande/bortfall och hjälpvariablerna. Det görs genom att beräkna skattad andel svarande i olika grupper, bestämda av respektive hjälpvariabel. Den totala svarsandelen har skattats till 65,0 procent.

Vid stora skillnader mellan svarsandelarna utgör variabeln en stark kandidat till hjälpvariabel.

Tabell 2. Skattad andel svarande fördelat på moderns födelse land

Moderns födelse land	Svarsandel (%)
Sverige	66,8
Övriga Världen	57,4
Totalt	65,0

Tabell 3 Skattad andel svarande fördelat på moderns ålder

Moderns ålder	Svarsandel (%)
28-39	55,3
40-45	67,1
46 -	70,8
Totalt	65,0

Tabell 4. Skattad andel svarande fördelat på region

Region	Svarsandel (%)
Storstadskommuner	63,4
Övriga kommuner	65,3
Totalt	65,0

Tabell 5. Skattad andel svarande fördelat på högsta utbildning föräldrar

Utbildningsnivå	Svarsandel (%)
Grundskoleutbildning	49,6
Gymnasial utbildning	58,3
Eftergymnasial utbildning	71,7
Totalt	65,0

Tabell 6. Skattad andel svarande fördelat på elevens kön

Kön	Svarsandel (%)
Pojke	65,0
Flicka	65,0
Totalt	65,0

Tabell 7. Skattad andel svarande fördelat på huvudman

Huvudman	Svarsandel (%)
Kommunal	65,0
Fristående	65,6
Totalt	65,0

Tabellerna 2-7 visar att variablerna födelseland, ålder och högsta utbildning föräldrar kan anses som starka (beträffande kriterium 1).

4.3 Kriterium 3: Variabeln avgränsar (viktiga) redovisningsgrupper

Om hjälpvariabeln avgränsar viktiga redovisningsgrupper kan kvaliteten bli bättre i dessa grupper. Framförallt blir skattningarna säkrare om hjälpvariabeln väl avgränsar redovisningsgruppen.

Variabeln elevens kön avgränsar redovisningsgrupper och bör därför vara med i hjälpvektorn om det är möjligt.

4.4 Slutligt val av hjälpvektor

Efter en sammanvägning av analysen kring de tre kriterierna samt efter kontroll av vikternas fördelning används följande hjälpvektor:

Elevens Kön + Moderns födelseland + Region + Högsta utbildning föräldrar

5 Teknisk beskrivning av urval och estimation

Vi har en population U bestående av N personer. De parametrar vi är intresserade av är främst funktioner av två totaler $Y = \sum_U y_k$ och $Z = \sum_U z_k$, där y_k är värdet på variabel y för person k och z_k värdet på en annan variabel för samma person. Vi kan definiera y (och även z) som en dikotom variabel, d.v.s.

$$y_k = \begin{cases} 1 & \text{om person } k \text{ har studerade egenskap} \\ 0 & \text{för övrigt} \end{cases} \quad (5.1)$$

Det finns givetvis också intresse av parametrar för olika redovisningsgrupper. Låt oss benämna dessa $U_1, \dots, U_d, \dots, U_D$, där

$U = \bigcup_{d=1}^D U_d$. Totalen för redovisningsgrupp d kan skrivas

$$Y_d = \sum_U y_{dk} \quad (5.2)$$

$$\text{där } y_{dk} = \begin{cases} y_k & \text{för } k \in U_d \\ 0 & \text{för övrigt.} \end{cases}$$

Z_d bildas på likartat sätt.

En generell parameter för redovisningsgrupp d (d kan också avse hela populationen) kan skrivas $\theta_d = C \frac{Y_d}{Z_d}$, där C är en konstant.

Den vanligaste parametern är en procentuell andel, som erhålles när $C = 100$ och $z_k = 1$ för alla k , och y är definierad enligt (5.1). Om vi låter N_d vara antalet personer i redovisningsgrupp d , då kan parametern skrivas

$$P_d = 100 \frac{\sum_U y_{dk}}{N_d} \quad (5.3)$$

Vi drar ett urval s av storleken n , men p.g.a. övertäckning och bortfall har vi endast svarsmängden r av storleken m att utföra beräkningarna på.

Den ”konventionella” estimatorn (för Y_d), har följande form:

$$\hat{Y}_d = \sum_r w_{1k} w_{2k} y_{dk} \quad (5.4)$$

där

$w_{1k} =$ totalt antal skolor i respektive stratum/
antal utvalda skolor i respektive stratum

$w_{2k} =$ totalt antal elever i respektive skola/
antal elever där föräldrarna svarat på enkäten i respektive skola

I estimator (5.4) används ingen ytterligare hjälpinformation än stratifieringsinformationen. Denna estimationsmetod brukar kallas ”rak uppräknning inom strata”.

I syfte att erhålla en estimator med mindre urvalsfel och bortfallsskevhet än estimator (5.4) utnyttjar vi hjälpinformation också i estimationen. Vi bildar en hjälpvektor \mathbf{x}_k , som anger till vilka kategorier av

Elevens Kön + Moderns födelseland + Region + Högsta utbildning föräldrar

som elev k tillhör. Från register framställer vi sedan hjälptotalerna $\sum_s \mathbf{x}_k$. Vi utnyttjar denna hjälpinformation i en kalibreringsestimator.

Kalibreringsestimatorn för totalen Y_d har följande utseende:

$$\hat{Y}_{wd} = \sum_r w_{1k} w_{2k} v_{2k} y_{dk} = \sum_r w_k y_{dk} \quad (5.5)$$

där

$$v_{2k} = 1 + \left(\sum_s \mathbf{x}_k - \sum_r w_{2k} \mathbf{x}_k \right)' \left(\sum_r w_{2k} \mathbf{x}_k \mathbf{x}_k' \right)^{-1} \mathbf{x}_k \quad (5.6)$$

Vid skattning av en parameter av typen $\theta_d = C \frac{Y_d}{Z_d}$ skattas respektive total

med hjälp av kalibreringsvikterna $w_k = w_{1k} w_{2k} v_{2k}$.

Referenser:

Lundström S. och Särndal C.-E. (2001). *Estimation in the Presence of Nonresponse and Frame Imperfection*. Stockholm: Statistics Sweden