

# Arbeidsnotater

S T A T I S T I S K S E N T R A L B Y R Å

Dronningensgt. 16, Oslo-Dep., Oslo l. Tlf. 41 38 20

IO 74/50

21. november 1974

## VARIANSESTIMERING FOR NIVÅTALLESTIMATER OG ENDRINGSTALLESTIMATER VED BYRÅETS ARBEIDSKRAFTUNDERSØKELSER

av

John Dagsvik<sup>\*)</sup>

	Side
1. Innledning .....	2
2. Estimater for nivåfall og endringstall, samt deres standard- avvik. Estimater for "design-effekten" og varianskomponentene mellom og innen utvalgsområdene .....	3
2.1. Estimatorene .....	3
2.2. Estimater for nivåfall og deres standardavvik .....	3
2.3. Varianskomponentene mellom og innen utvalgsområdene .....	5
2.4. Estimater for endringstall og deres standardavvik .....	6
2.5. Estimater for "design-effekten" .....	7
3. Den matematiske teori for estimering av variansen til endrings- tallestimater når utvalget er roterende .....	8
3.1. Varianser .....	8
3.2. Estimatorer .....	13
3.3. Spesialtilfellet Oslo .....	16
3.4. Omforminger og tilnærmelser i beregningsprogrammet .....	17
Referanser .....	18

\*) Jeg vil takke Ib Thomsen, Petter Laake, Stein Østerlund Petersen og medarbeiderne i Sosiodemografisk forskningsgruppe for nyttige kommentarer under bearbeidelsen av manuskriptet.

Kjetil Sørliie er ansvarlig for programmeringen av estimatorene, og jeg vil derfor rette en spesiell takk til han for en meget effektiv hjelp og et godt samarbeid.

*Ikke for offentliggjøring. Dette notat er et arbeidsdokument og kan siteres eller refereres bare etter spesiell tillatelse i hvert enkelt tilfelle. Synspunkter og konklusjoner kan ikke uten videre tas som uttrykk for Statistisk Sentralbyrås oppfatning.*

## Sammendrag

Ved Byråets arbeidskraftundersøkelser (AKU) skiftes en del av utvalget ut hvert kvartal i følge en spesiell rotasjonsplan. I dette notatet har en med utgangspunkt i denne rotasjonsplanen, utledet en estimator for variansen til estimert endring i tallene fra en undersøkelse til en annen.

Ved Intervjukontoret ble det våren 1974 laget et program for estimering av varianser til estimatorer for nivå-tall. Dette programmet er utvidet med en subrutine som beregner tilsvarende varians-estimatorer for endringstall basert på den estimatoren vi har utledet.

Det utvidede variansberegningsprogrammet er anvendt til beregning av estimatorer for nivå-tall, endringstall samt deres varianser. Videre har en funnet estimatorer for "design-effekten" og for variansenkomponentene som skyldes variansen mellom utvalgsområdene og innen utvalgsområdene.

## 1. Innledning

Mens variansestimatorene ved utvalgsundersøkelser der utvalget er trukket rent tilfeldig er forholdsvis enkle og kan beregnes for hånd, er dette en temmelig arbeidskrevende regneoperasjon når utvalget er trukket etter Byråets utvalgsplan.

Variansestimering på grunnlag av slike kompliserte utvalgsplaner har ofte vært basert på grove tilnæringsformler en har funnet blant annet ved empiriske undersøkelser i andre land.

Ved Intervjukontoret ble det våren 1974 laget et variansberegningsprogram på grunnlag av en forventningsrett variansestimator for to-trinns utvalg. (Laake et al. 1974). Dette programmet estimerer variansen til estimert gjennomsnitt og samlet antall av en variabel, ved en enkelt undersøkelse. Den formelle utledningen av denne variansestimatorene er gjengitt blant annet hos Des Raj 1968 og hos Hoem 1973. I AKU benyttes en roterende utvalgsplan som er slik at en fjerdedel av utvalget skiftes ut ved hvert kvartal og hvor halvparten av et ordinært kvartalvis utvalg er felles ved ett kvartal ett år og samme kvartal neste år. Intervjukontorets opprinnelige variansberegningsprogram kan ikke uten videre benyttes til å beregne variansestimater til estimatene for endringene i tallene fra en undersøkelse til en annen.

Vi har imidlertid utledet en variansestimator basert på en slik rotasjonsplan. Denne estimatoren har en slik form at de numeriske beregningene kan gjøres ved hjelp av en subrutine til det opprinnelige variansberegningsprogram.

I Byrådet er det tidligere også utledet en slik variansestimator men denne er basert på grove tilnæringer (Østerlund Petersen 1972).

I den første delen av notatet presenteres standardavvikestimater for nivåtallestimater og endringstallestimater. Videre gis estimater for "designeffekten" og varianskomponentene mellom og innen utvalgsområdene. I den andre delen presenteres den formelle utledningen av estimatoren for variansen til endringstallestimatoren, dvs. differensen mellom de respektive nivåtallsestimatorene.

## 2. Estimater for nivåfall og endringstall, samt deres standardavvik. Estimater for "design-effekten" og varianskomponentene mellom og innen utvalgsområdene

### 2.1. Estimatorene

De variansestimatorene som er programmert gjelder egentlig bare når nivåfallene estimeres ved vanlig oppblåsning av tallene i utvalget. I Byråets AKU er imidlertid estimeringsmetoden en annen idet en blåser opp tallene i utvalget for hvert kjønn og ettårige aldersklasser og deretter summerer disse. Nå gir variansberegningsprogrammet også utskrift av nivåfallene estimert ved vanlig oppblåsning av tallene fra utvalget slik at det er mulig å sammenlikne de to estimeringsmetodene (og eventuelt avgjøre om resultatene er signifikant forskjellige). Det følger forøvrig av Dagsvik 1974 (setning 3.7 (iv)) at variansen til Byråets AKU-estimater er større enn variansen til estimater beregnet ved vanlig oppblåsning.

Variansberegningsprogrammet gir endelig estimater for "design-effekten". Denne er definert som forholdet mellom variansen til den estimator som benyttes og variansen til den tilsvarende estimator dersom utvalget var trukket enkelt tilfeldig.

I Byråets AKU er annentrinns trekkeenheter de enkelte husholdninger. I den matematiske teorien inngår antall husholdninger i de uttrukne utvalgsområdene i beregningen av den totale utvalgsbrøk som benyttes i variansestimatorene og ved oppblåsingen. Imidlertid inneholder ikke intervjukontorets utvalgsregister husholdning som enhet slik at vi må nøye oss med å estimere den totale utvalgsbrøk på grunnlag av kjennskapet til antall personer i bestanden. Vi innfører dermed en ny feilkilde i tillegg til de mer tradisjonelle feilkilder som frafall, utvalgsfeil og registerfeil. Vi har imidlertid antatt at denne feilkilden har liten betydning sammenliknet med de andre feilkildene.

### 2.2. Estimater for nivåfall og deres standardavvik

Nivåfallene ved første og annet kvartal 1973 og 1974 for sysselsatte i hovednæringer estimert ved de to omtalte metodene samt standardavvikestimater for oppblåste tall er gjengitt i tabell 1. Denne viser at forskjellen mellom estimatene beregnet ved de to estimeringsmetodene er liten. Forskjellen er størst for jordbruk, men heller ikke her er forskjellen større enn ett standardavvik. Dersom standardavvik brukes som signifikansmål (Thomsen 1973) finner vi at resultatene fra de to

Tabell 1. Sysselsatte etter næring. 1 000

Næringer	1. kvartal 1973			2. kvartal 1973			1. kvartal 1974			2. kvartal 1974		
	AKU tall	Opp- blåste tall	Stand- ard avvik									
Jordbruk .....	150	156	10,5	163	167	11,5	121	129	9,0	142	151	10,2
Skogbruk .....	19	19	3,0	14	15	2,7	14	14	2,6	13	14	2,4
Fiske og fangst .....	31	31	4,2	21	22	3,4	35	35	4,8	24	23	3,6
Bergverksdrift .....	10	10	3,7	14	14	4,7	12	11	3,0	8	8	2,9
Industri .....	378	378	14,5	390	391	15,5	378	378	14,4	399	400	14,6
Kraft og vannfor- syning .....	13	13	2,0	16	16	2,3	17	18	2,3	16	17	2,4
Bygg og anlegg .....	135	135	7,3	142	143	7,4	140	140	6,7	149	150	7,3
Varehandel, hotell og restaurant .....	269	270	11,9	279	276	11,8	281	282	11,7	274	273	11,8
Transport, lagring, post og telekomm. ...	159	157	7,4	161	160	8,6	159	156	7,6	164	162	8,2
Bank og finansierings- virksomhet, eiendoms- drift og forr. tjenesteyting .....	67	66	4,8	63	60	4,8	67	66	5,8	66	64	5,3
Offentlig, sosial og privat tjenesteyting.	411	411	15,1	393	392	14,4	420	417	15,2	391	388	14,7
Sysselsatte i alt, medregnet uoppgift ..	1 645	1 649	24,1	1 657	1 663	24,3	1 645	1 652	24,5	1 649	1 643	25,0

estimeringsmetodene ikke er signifikant forskjellige. Studerer vi standardavvikestimaterne ser vi at de varierer forholdsvis lite over tid for samme næring mens de varierer en del mellom enkelte næringer hvor antall sysselsatte er av samme størrelsesorden.

### 2.3. Varianskomponentene mellom og innen utvalgsområdene

Variansen kan oppdeles i komponenter som måler henholdsvis variasjonen innen og mellom utvalgsområdene. Varianskomponenten mellom utvalgsområdene er derfor et mål på hvor god stratifiseringen er. Dessverre gir ikke det opprinnelige variansberegningssystemet utskrift av disse varianskomponentene for nivåtallestimatene. Subrutinen som beregner variansestimaterne for endringstallestimatene gir derimot utskrift av tilsvarende komponenter. Disse er gjengitt i tabell 2. Tallene tyder på at variansen mellom utvalgsområdene er av samme størrelsesorden som variansen innen utvalgsområdene for jordbruk, industri, varehandel m.m., og offentlig sosial og privat tjenesteyting, mens variansen mellom utvalgsområdene er vesentlig mindre enn variansen innen utvalgsområdene for kraft og vannforsyning, og bygg og anlegg. I Sverige er det gjort tilsvarende beregninger og disse resultatene viser gjennomgående at varianskomponenten innen utvalgsområdene er vesentlig større enn mellom utvalgsområdene for de fleste næringer. Dette tyder på at stratifiseringen ved den svenske utvalgsplanen er bedre enn ved den norske.

Tabell 2. Varianskomponentene innen- og mellom utvalgsområdene for 2. kvartal 1973 og 1974 i prosent av total varians

Næringer		Variansen innen utvalgsområdene	Variansen mellom utvalgsområdene
Jordbruk	1973 .....	50	50
	1974 .....	56	44
Industri	1973 .....	51	49
	1974 .....	53	47
Kraft og vannforsyning	1973 .....	64	36
	1974 .....	74	26
Bygg og anlegg	1973 .....	76	24
	1974 .....	81	19
Varehandel, hotell og restaurant	1973 .....	54	46
	1974 .....	52	48
Off. sosial og privat tjenesteyting	1973 .....	57	43
	1974 .....	52	48

#### 2.4. Estimater for endringstall og deres standardavvik

Tabell 3 viser estimerte endringer fra første og annet kvartal 1973 til henholdsvis første og annet kvartal 1974 for noen hovednæringer, samt estimert standardavvik for differensen mellom de respektive oppblåste tall. Vi legger merke til at bortsett fra jordbruk er endringstallestimatene i annet kvartal mindre enn standardavvikestimaterne. Altså er endringstallestimatene i annet kvartal signifikante bare for jordbruket.

Videre er standardavvikestimaterne for endringstallestimatene stort sett mindre enn standardavvikestimaterne for de respektive nivå-tallestimater. I samsvar med resultatene i kapitel 2.2 viser også standardavvikestimaterne for endringstallestimatene stor stabilitet over tid.

I Byrådet er det også gjort sammenlikninger med en annen metode for estimering av endringstall, nemlig ved å estimere endringene på grunnlag av den delen av utvalget som er felles ved begge undersøkelser. Siden denne delen av utvalget bare utgjør halvparten av et ordinært AKU-utvalg har følgelig estimatene større varians. En har kalt denne metoden den direkte estimeringsmetoden mens den andre metoden vi har brukt er blitt kalt den indirekte metoden. For endringene fra 1. kvartal 1973 til 1. kvartal 1974 gir de to metodene en nedgang på henholdsvis 27 000 og 15 000 for jordbruket. Det første estimatet har et standardavvik på ca. 10 000 og det andre estimatet har et standardavvik større enn 10 000. Begge konfidensintervallene for de respektive estimatene inneholder altså intervallet (17 000, 25 000) slik at vi ikke kan si at de to metodene gir signifikant forskjellig resultat i dette tilfelle. Tilsvarende betraktninger gjelder også for andre næringer.

Den store forskjellen mellom estimatene ved de to estimeringsmetodene gir likevel grunn til å undersøke om en kombinasjon av indirekte og direkte estimater kan gi bedre resultater. Slike "sammensatte estimater" er i bruk i andre land og har vist seg å gi presisjonsgevinster for mange variable.

Tabell 3. Endringstall 1973-1974. 1 000

Næringer	Endringer 1. kvartal 1973 til 1. kvartal 1974			Endringer 2. kvartal 1973 til 2. kvartal 1974		
	AKU tall	Opp- blåste tall	Standard- avvik	AKU tall	Opp- blåste tall	Standard- avvik
	Jordbruk .....	-29	-27	9,7	-21	-16
Industri .....	0	0	12,7	9	9	13,6
Kraft og vannforsyning ..	4	5	2,6	0	-1	3,0
Bygg og anlegg .....	5	5	7,0	7	7	8,8
Varehandel, hotell og restaurant .....	12	12	11,7	-5	-3	7,0
Off. sosial og privat tjenesteyting .....	9	6	13,6	-2	-4	12,2

### 2.5. Estimerer for "design-effekten"

"Design-effekten" er et mål for "effekten" av den spesielle utvalgsplanen som ligger til grunn relativt til en utvalgsplan der utvalget er trukket reñt tilfeldig. Dersom et kjennetegn er jevnt fordelt i bestanden er det av mindre betydning hvilken utvalgsplan som ligger til grunn. Derfor er "design-effekten" også et mål på hvor homogen bestanden er med hensyn på de respektive kjennetegn.

Ved å studere tabell 4 ser vi at tre næringsgrupper skiller seg klart ut med stor "design-effekt", nemlig gruppene jordbruk, fiske og fangst, og bergverksdrift. Spesielt bergverksdrift har en svært stor "design-effekt". Dette er i samsvar med det kjennskap en har fra andre kilder om næringens geografiske fordeling. (Rideng 1970). Bortsett fra bergverksdrift og fiske og fangst viser disse estimatene stabilitet over tid.

For de 5 siste næringsgruppene i tabellen ligger estimatene stort sett i underkant av 1,5.

I Byrået har en ofte brukt en variansestimator basert på antagelsen at design-effekten er lik 1,5. Våre resultater viser at denne estimatoren ville ha underestimert variansen for de 4 første næringsgruppene og overestimert variansen for de resterende gruppene.

Tabell 4. Designeffekt

Næringer	1. kvar- tal 1973	2. kvar- tal 1973	1. kvar- tal 1974	2. kvar- tal 1974
Jordbruk .....	2,6	2,6	2,4	2,6
Skogbruk .....	1,7	1,8	1,7	1,6
Fiske og fangst .....	2,0	2,0	2,5	2,1
Bergverksdrift .....	4,8	5,8	3,1	4,2
Industri .....	1,6	1,8	1,5	1,5
Kraft og vannforsyning .....	1,2	1,3	1,1	1,2
Bygg og anlegg .....	1,3	1,4	1,2	1,2
Varehandel, hotell og restaurant .	1,5	1,5	1,3	1,3
Transport, lagring, post og telekomm. ....	1,3	1,5	1,3	1,3
Bank og finansieringsvirksomhet, eiendomsdrift og forr. tjeneste- yting .....	1,3	1,4	1,7	1,5
Off. sosial og privat tjeneste- yting .....	1,5	1,4	1,2	1,3

### 3. Den matematiske teori for estimering av variansen til endringstall- estimerer når utvalget er roterende

#### 3.1. Varianser

Utgangspunktet for utledningen av en estimator for variansen til endringstallestimatoren er at det samlede utvalg som var med i en av undersøkelsene på tidspunkt 1 og tidspunkt 2 kan betraktes som ett utvalg trukket uten tilbakelegging ved første undersøkelsestidspunkt. Dette utvalget består derfor av 3 delutvalg,  $S_1$  som bare er med i første undersøkelse,  $S_2$  som er med i begge undersøkelser, og  $S_3$  som er med i siste undersøkelse. I samsvar med gjeldende utvalgsplan antar vi at delutvalget  $S_1US_2$  er like stort som delutvalget  $S_2US_3$ , nemlig et ordinært AKU utvalg på  $n$  trekkenheter. Størrelsen på delutvalget  $S_2$  betegner vi med  $n'$ . I det følgende er symboler og konvensjoner definert i samsvar med Hoem 1973.

Bestanden er inndelt i et visst antall strata hvor stratum nr.  $i$  inneholder  $M_i$  primære utvalgsområder. Det  $j$ -te primære utvalgsområde

i stratum nr.  $i$  har  $N_i(j)$  trekkeenheter. Den  $k$ -te trekkeenheter i  $j$ -te utvalgsområde i  $i$ -te stratum har verdien  $a_i(j, k)$  ved tidspunkt 1 og verdien  $c_i(j, k)$  ved tidspunkt 2. Vi antar at det er like mange trekkeenheter i de primære utvalgsområdene ved begge de to undersøkelsestidspunktene. Vi innfører notasjonene

$$N_i = \sum_j N_i(j), N = \sum_j N_i,$$

$$a_i(j) = \sum_k a_i(j, k), \bar{a}_i(j) = a_i(j) / N_i(j),$$

$$a_i = \sum_j a_i(j), \bar{a}_i = a_i / M_i,$$

$$a = \sum_i a_i,$$

og tilsvarende notasjoner for  $c$ -verdiene.

Blant de  $M_i$  utvalgsområdene i stratum  $i$  trekkes  $m_i$  områder med nummerene  $J_{i1}, J_{i2}, \dots, J_{im}$ , rent lotterisk. Sannsynligheten for at et vilkårlig område i  $i$ -te stratum skal komme med i utvalget er altså  $m_i/M_i = \pi_i$ . Vi lar  $J$  være vektoren med komponenter lik numrene på alle utvalgsområdene som trekkes. Disse områdene er de samme ved begge undersøkelsestidspunktene. På tidspunkt 1 trekkes  $n_{ir}(J)$  trekkeenheter rent lotterisk fra utvalgsområde  $J_{ir}$  hvor  $S_1$  består av de  $n_{ir}(J) - n_{ir}'(J)$  første,  $S_2$  består av de  $n_{ir}'(J)$  neste og  $S_3$  består av de  $n_{ir}(J) - n_{ir}'(J)$  resterende trekkeenheter,  $r = 1, 2, \dots, m_i$ ,  $i = 1, 2, \dots$ . Numrene på de enhetene som trekkes ut fra område  $J_{ir}$  betegnes  $K_{ir1}, K_{ir2}, \dots$ , og vi definerer indeksvariablene

$$U_{irs} = \begin{cases} 1 & \text{dersom } s \leq n_{ir}(J) \\ 0 & \text{ellers} \end{cases}$$

og

$$V_{irs} = \begin{cases} 1 & \text{dersom } s > n_{ir}(J) - n_{ir}'(J) \\ 0 & \text{ellers} \end{cases}$$

Videre innfører vi

$$X_{irs} = a_i (J_{ir}, K_{irs}), Y_{irs} = c_i (J_{ir}, K_{irs}),$$

$$X_{ir} = \sum_s U_{irs} X_{irs}, Y_{ir} = \sum_s V_{irs} Y_{irs},$$

$$\bar{X}_{ir} = X_{ir}/n_{ir} (J), \bar{Y}_{ir} = Y_{ir}/n_{ir} (J),$$

$$X_i = \sum_r X_{ir}, Y_i = \sum_r Y_{ir},$$

$$\bar{X}_i = X_i/m_i, \bar{Y}_i = Y_i/m_i,$$

$$X = \sum_i X_i, Y = \sum_i Y_i,$$

$$b_i (J) = \pi_i b_i (J).$$

$b_i (J)$  er utvalgsbrøken for stratum  $i$  og er den samme for alle utvalgsområdene innen stratumet, og  $b (J)$  er den totale utvalgsbrøk og er uavhengig av  $i$ . De estimatene vi tar utgangspunkt i er

$$\hat{a}_i = b^{-1} (J) X_i, \hat{c}_i = b^{-1} (J) Y_i$$

$$\hat{a} = \sum_i \hat{a}_i \text{ og } \hat{c} = \sum_i \hat{c}_i,$$

som er estimatorer for henholdsvis  $a_i$ ,  $c_i$ ,  $a$  og  $c$ . Estimatorene  $\hat{a} - \hat{c}$  er derfor en estimator for endringen  $a - c$ , av samlet variabelverdi. Vi skal utlede variansen til  $\hat{a} - \hat{c}$  og angi en estimator for denne variansen. Siden  $\text{Var} \{\hat{a} - \hat{c}\} = \text{Var} \hat{a} + \text{Var} \hat{c} - 2 \text{cov} \{\hat{a}, \hat{c}\}$  og de to første leddene er variansene til nivåtallestimatorene på de to tidspunktene (som er kjente), er det nok å finne  $\text{cov} \{\hat{a}, \hat{c}\}$ . Denne kovariansen kan splittes opp på samme måte som variansen til nivåtallestimatorene, nemlig i komponenter som er et uttrykk for kovariansen henholdsvis innen og mellom utvalgsområdene.

Setning 1: Kovariansen mellom utvalgsområdene er gitt ved

$$\text{cov} \{E(\hat{a}|J), E(\hat{c}|J)\} = \sum_i \pi_i^{-1} (M_i - m_i) \xi_i$$

hvor

$$\xi_i = \sum_j (a_i(j) - \bar{a}_i)(c_i(j) - \bar{c}_i) / (M_i - 1).$$

Bevis: Vi har at

$$E \{ \hat{a}_i | J = j \} = \sum_r \pi_i^{-1} a_i(j_{ir})$$

og

$$E \{ \hat{c}_i | J = j \} = \sum_r \pi_i^{-1} c_i(j_{ir}).$$

Ved å anvende Sverdrup 1963 (side 328) får vi umiddelbart at

$$\text{cov} \{E(\hat{a}_i | J), E(\hat{c}_i | J)\} = \pi_i^{-1} (M_i - m_i) \xi_i.$$

Samme resonnement som Hoem 1973 (side 12) gir videre

$$\text{cov} \{E(\hat{a}_i | J), E(\hat{c}_j | J)\} = 0 \text{ når } i \neq j.$$

Dermed er setningen bevist.  $\square$

Merknad: Det fremgår av setningen at kovariansen mellom utvalgsområdene er uavhengig av størrelsen på delutvalget  $S_2$ .

Lemma 2:

$$\begin{aligned} \text{(i)} \quad & \sum_s E \{ U_{irs} V_{irs} X_{irs} Y_{irs} | J = j \} \\ & = n_{ir}^1(j) (N_i(j_{ir}) - 1) \xi_i(j_{ir}) / N_i(j_{ir}) + n_{ir}^1(j) \bar{a}_i(j_{ir}) \bar{c}_i(j_{ir}) \end{aligned}$$

$$\begin{aligned} \text{(ii)} \quad & E \{ X_{ir} Y_{ir} | J = j \} = n_{ir}^2(j) \bar{a}_i(j_{ir}) \bar{c}_i(j_{ir}) \\ & + \xi_i(j_{ir}) \{ n_{ir}^1(j) - n_{ir}^2(j) / N_i(j_{ir}) \} \end{aligned}$$

hvor

$$\xi_i(j) = \sum_k (a_i(j,k) - \bar{a}_i(j))(c_i(j,k) - \bar{c}_i(j)) / (N_i(j) - 1).$$

Bevis: For  $J = j$  har hver trekkeenhet innen utvalgsområde nr.  $j_{ir}$  sannsynligheten  $1/N_i(j_{ir})$  for å bli trukket ut i en trekning. Herav følger at

$$\begin{aligned} E \{X_{irs} Y_{irs} | J = j\} &= \sum_k a_i(j_{ir}, k) c_i(j_{ir}, k) \Pr\{K_{irs} = k | J = j\} \\ &= \sum_k a_i(j_{ir}, k) c_i(j_{ir}, k) / N_i(j_{ir}). \end{aligned}$$

Siden  $S_2$  inneholder  $n_{ir}^!(j)$  enheter fra utvalgsområde nr.  $j_{ir}$  er (i) dermed bevist.  $\square$

(ii): For  $J = j$  har to ulike enheter innen utvalgsområde nr.  $j_{ir}$  sannsynligheten  $1/N_i(j_{ir}) (N_i(j_{ir}) - 1)$  for å bli trukket ut i en trekning. Siden  $U_{irs} V_{irp}$  antar verdien 1 for  $n_{ir}^2(j) - n_{ir}^!(j)$  verdier av  $(s, p)$  når  $s \neq p$ , blir

$$\begin{aligned} \sum_{s \neq p} E \{U_{irs} V_{irp} X_{irs} Y_{irp} | J = j\} \\ = \{n_{ir}^2(j) - n_{ir}^!(j)\} \sum_{k \neq h} a_i(j_{ir}, k) c_i(j_{ir}, h) / N_i(j_{ir}) (N_i(j_{ir}) - 1) \end{aligned}$$

Sammen med (i) gir dette det ønskede resultat.  $\square$

Det neste resultatet følger nå umiddelbart fra lemma 2.

Setning 3: Kovariansen innen utvalgsområdene er gitt ved

$$E \text{ cov} \{\hat{a}, \hat{c} | J\} =$$

$$E \{b^{-2}(J) \sum_{i,r} [n_{ir}^!(J) - n_{ir}^2(J) / N_i(J_{ir})] \xi_i(J_{ir})\}$$

Av setning 3 følger umiddelbart:

Korollar 4: Dersom  $n_{ir}^!(J) = n_{ir}(J)/q$  hvor  $q$  er uavhengig av  $i$  og  $r$  for alle  $i$  og  $r$  er

$$E \text{ cov} \{\hat{a}, \hat{c} | J\} =$$

$$E \{b^{-1}(J) \sum_{i,r} \pi_i^{-1} N_i(J_{ir}) (q^{-1} - b_i(J)) \xi_i(J_{ir})\}.$$

I Byråets rotasjonsplan er halvparten av et ordinært AKU utvalg felles ved en undersøkelse et kvartal og undersøkelsen samme kvartal året etter.

Ved å sette  $q = 2$  får vi derfor kovariansen innen utvalgsområdene i dette tilfelle.

Merknad: Vi ser av setning 3 at dersom delutvalget  $S_2$  blir lite nok skifter kovariansen innen utvalgsområdene fortegn.

Korollar 5: Når  $S_2$  er tom er

$$\text{Cov} \{ \hat{a}, \hat{c} \} = \sum_i \left[ \pi_i^{-1} (M_i - m_i) \xi_i - \pi_i^{-1} \sum_j N_i(j) \xi_i(j) \right].$$

Forholdet mellom det negative leddet i korollar 5 og kovariansen innen utvalgsområdene er av størrelsesorden  $2\pi_i^{-1} b(J) \approx 12b(J)$  når  $q = 2$ .

I mange tilfeller kan vi derfor regne kovariansen tilnærmet lik kovarianskomponenten mellom utvalgsområdene når  $S_2$  er tom.

### 3.2. Estimatorer

Vi skal nå behandle estimeringen av kovariansen mellom  $\hat{a}$  og  $\hat{c}$  i det tilfelle at delutvalget  $S_2$  er ikke-tomt. De estimatorene vi kommer fram til kan imidlertid ikke brukes når  $S_2$  er tom.

For å vise at estimatorene er forventningsrette trenger vi følgende hjelperesultater:

Lemma 6: Estimatoren

$$W_i = b^{-2} \sum_r \frac{\{n_{ir}^1(J) - n_{ir}^2(J)/N_i(J_{ir})\}}{n_{ir}^1(J) \{1 - n_{ir}^1(J)/n_{ir}^2(J)\}} \sum_s U_{irs} V_{irs} (X_{irs} - \bar{X}_{ir})(Y_{irs} - \bar{Y}_{ir})$$

er forventningrett for

$$E \{ b^{-2} \sum_r \left[ n_{ir}^1(J) - n_{ir}^2(J) / N_i(J_{ir}) \right] \xi_i(J_{ir}) \}.$$

Bevis: Ved å bruke lemma 2 får vi

$$\begin{aligned}
 & E \left\{ \sum_s U_{irs} V_{irs} (X_{irs} - \bar{X}_{ir})(Y_{irs} - \bar{Y}_{ir}) \mid \mathcal{J} = j \right\} \\
 &= E \left\{ \sum_s U_{irs} V_{irs} X_{irs} Y_{irs} - n'_{ir}(j) \bar{X}_{ir} \bar{Y}_{ir} \mid \mathcal{J} = j \right\} \\
 &= \xi_i(j_{ir}) n'_{ir}(j) \{1 - n'_{ir}(j)/n_{ir}^2(j)\} \quad \square .
 \end{aligned}$$

Lemma 7:

$$\begin{aligned}
 & E \left\{ b^{-2}(j) \sum_r (X_{ir} - \bar{X}_i)(Y_{ir} - \bar{Y}_i) \right\} m_i / (m_i - 1) \\
 &= \sum_r E \left\{ b^{-2}(j) \left[ n'_{ir}(j) - n_{ir}^2(j) / N_i(j_{ir}) \right] \xi_i(j_{ir}) \right\} \\
 &+ M_i \pi_i^{-1} \xi_i .
 \end{aligned}$$

Bevis: Lemma 2 gir oss

$$\begin{aligned}
 & \sum_r \text{cov} \{ b^{-1}(j) X_{ir}, b^{-1}(j) Y_{ir} \} = \\
 & \sum_r E \{ X_{ir} b^{-1}(j), Y_{ir} b^{-1}(j) \} - a_i c_i / m_i \\
 &= \sum_r E \left\{ b^{-2}(j) \left[ n'_{ir}(j) - n_{ir}^2(j) / N_i(j_{ir}) \right] \xi_i(j_{ir}) \right\} \\
 &+ \pi_i^{-1} (M_i - 1) \xi_i .
 \end{aligned}$$

Sammen med setning 1 og setning 3 får vi derfor

$$\begin{aligned}
 & \sum_{r,i} E \{ (X_{ir} - \bar{X}_i)(Y_{ir} - \bar{Y}_i) b^{-2}(j) \} m_i \\
 &= \sum_{i,r} \text{cov} \{ X_{ir} b^{-1}(j), Y_{ir} b^{-1}(j) \} m_i
 \end{aligned}$$

$$\begin{aligned}
& - \sum_i \text{cov} \{X_i b^{-1}(J), Y_i b^{-1}(J)\} \\
& = \sum_i [(m_i - 1) \sum_r E \{b^{-2}(J) [n'_{ir}(J) - n^2_{ir}(J)/N_i(J_{ir})] \xi_i(J_{ir})\} \\
& + M_i(M_i - 1) \xi_i - \pi_i^{-1} (M_i - m_i) \xi_i].
\end{aligned}$$

Da de to siste leddene reduserer seg til  $M_i \pi_i^{-1} (m_i - 1) \xi_i$ , er lemmaet bevist.

□.

Setning 8: Estimatoren

$$\begin{aligned}
H & = b^{-2}(J) \sum_i (1 - \pi_i) m_i \sum_r (X_{ir} - \bar{X}_i)(Y_{ir} - \bar{Y}_i) / (m_i - 1) \\
& + \sum_i \pi_i W_i
\end{aligned}$$

er forventningsrett for  $\text{cov} \{\hat{a}, \hat{c}\}$ .

Bevis: Av lemma 6 og lemma 7 følger det at

$$\begin{aligned}
E \{H \mid J = j\} & = E \sum_i [\pi_i b^{-2}(J) \sum_r \{n'_{ir}(J) - n^2_{ir}(J)/N_i(J_{ir})\} \xi_i(J_{ir}) \\
& + (1 - \pi_i) b^{-2}(J) \sum_r \{n'_{ir}(J) - n^2_{ir}(J)/N_i(J_{ir})\} \xi_i(J_{ir}) \\
& + (1 - \pi_i) M_i \pi_i^{-1} \xi_i] = E \text{cov} \{\hat{a}, \hat{c} \mid J\} \\
& + \text{cov} \{E(\hat{a} \mid J), E(\hat{c} \mid J)\} = \text{cov} \{\hat{a}, \hat{c}\} \quad \square.
\end{aligned}$$

Korollar 9: Dersom  $n'_{ir}(J) = n_{ir}(J)/q$  for alle  $i$  og  $r$  reduseres estimatoren  $H$  til

$$\begin{aligned}
H & = b^{-2}(J) \sum_i (1 - m_i/M_i) m_i \sum_r (X_{ir} - \bar{X}_i)(Y_{ir} - \bar{Y}_i) / (m_i - 1) \\
& + qb^{-1}(J) \sum_i \{1 - qb_i(J)\} \sum_{r,s} N_i(J_{ir}) U_{irs} V_{irs} (X_{irs} - \bar{X}_{ir})(Y_{irs} - \bar{Y}_{ir}) / (qn_{ir}(J) - 1).
\end{aligned}$$

### 3.3. Spesialtilfellet Oslo

I Oslo trekkes et utvalg på  $2n_o(J) - n_o'(J)$  trekkeenheter rent lotterisk fra en bestand på  $N_o$  trekkeenheter. Vi lar

$$a_o = \sum_k a_o(k) \text{ og } c_o = \sum_k c_o(k)$$

være samlet variabelverdi ved henholdsvis første og andre undersøkelsestidspunkt. Analogt til notasjonene i avsnittet foran lar vi  $K_{o1}, K_{o2}, \dots$ , være numrene på de enhetene som trekkes ut og vi innfører videre notasjonene

$$U_{os} = \begin{cases} 1 & \text{når } s \leq n_o(J) \\ 0 & \text{ellers} \end{cases}$$

$$V_{os} = \begin{cases} 1 & \text{når } s > n_o(J) - n_o'(J) \\ 0 & \text{ellers} \end{cases},$$

$$X_{os} = a_o(K_{os}), Y_{os} = c_o(K_{os}),$$

$$X_o = \sum_s U_{os} X_{os}, Y_o = \sum_s V_{os} Y_{os},$$

$$\bar{X}_o = X_o/n_o(J) \text{ og } \bar{Y}_o = Y_o/n_o(J).$$

Størrelsen på utvalget er valgt slik at

$$n_o(J) = b(J)N_o.$$

Vår estimator for endringstallene innen Oslo er  $\hat{a}_o - \hat{c}_o$  hvor

$$\hat{a}_o = b^{-1}(J) X_o \text{ og } \hat{c}_o = b^{-1}(J) Y_o.$$

Tilsvarende beviset for setning 3 og lemma 6 finner vi

#### Setning 10:

$$\text{cov} \{ \hat{a}_o, \hat{c}_o \} = E \{ b^{-2}(J) (n_o'(J) - n_o^2(J)/N_o) \} \xi_o$$

hvor

$$\xi_o = \sum_k (a_o(k) - \bar{a}_o)(c_o(k) - \bar{c}_o)/(N_o - 1).$$

Setning 11: Kovariansen mellom  $\hat{a}_0$  og  $\hat{c}_0$  estimeres ved

$$W_0 = b^{-2} (J) \frac{(n'_0(J) - n^2(J)/N_0)}{n'_0(J) \{1 - n'_0(J)/n^2(J)\}} \sum_s U_{os} V_{os} (X_{os} - \bar{X}_0)(Y_{os} - \bar{Y}_0).$$

Korollar 12: Dersom  $n'_0(J) = n_0(J)/q$  reduseres  $W_0$  til

$$W_0 = qb^{-1} (J) \{1 - qb(J)\} N_0 \sum_s U_{os} V_{os} (X_{os} - \bar{X}_0)(Y_{os} - \bar{Y}_0) / (qn_0(J) - 1).$$

### 3.4. Omforminger og tilnærmelser i beregningsprogrammet

Under programmeringen har en brukt tilnærmelsen

$$2n_{ir}(J) / (2n_{ir}(J) - 1) \approx 1 \text{ for alle } i \text{ og } r, \text{ og } 2n_0 / (2n_0 - 1) \approx 1.$$

Dermed forenkles estimatoren H til

$$H = b^{-2} (J) \sum_i (1 - \pi_i) m_i \sum_r (X_{ir} - \bar{X}_i)(Y_{ir} - \bar{Y}_i) / (m_i - 1) \\ + \sum_i \pi_i^{-1} \{1 - 2b_i(J)\} \sum_{r,s} U_{irs} V_{irs} (X_{irs} - \bar{X}_{ir})(Y_{irs} - \bar{Y}_{ir}).$$

Videre har en benyttet at

$$\sum_j 2(X_j - \bar{X})(Y_j - \bar{Y}) \\ = \sum_j \{(X_j + Y_j - \bar{X} - \bar{Y})^2 - (X_j - \bar{X})^2 - (Y_j - \bar{Y})^2\}.$$

Referanser

- Dagsvik, J. 1974: Etterhåndsstratifisering og estimering innen delbestander. Statistisk Sentralbyrå, Artikkel 66.
- Des Raj 1968: Sampling theory. McCraw-Hill, London.
- Hoem, J. M. 1973: Statistisk Sentralbyrås utvalgsundersøkelser: Elementer av det matematiske grunnlaget. Statistisk Sentralbyrå, Artikkel 58.
- Laake, P. Forsén, L. og Sørli, K. 1974: Orientering om bruk av Byråets program for estimering av gjennomsnitt og estimering av variansen til estimatoren i Byråets intervjuundersøkelser. Stensil. (PL/TF, 24/4-1974).
- Rideng, A. 1970: Typifisering av kommuner i Norge, på grunnlag av næringsstrukturen 1968. Meddelelser fra Geografisk Institutt, Universitetet i Oslo. Kulturgeografisk serie, nr. 2.
- Sverdrup, E. 1964: Lov og tilfeldighet. Bind I. Universitetsforlaget, Oslo.
- Thomsen, I. 1973: Hvor oppdelt kan en offentliggjøre resultatene fra en intervjuundersøkelse? Statistisk Sentralbyrå, Metodehefte nr. 2 (IO 73/6).
- Østerlund Petersen, S. 1972: Arbeidskraftundersøkelsene. Om endringer i tallene fra en undersøkelse til en annen. Statistisk Sentralbyrå, Metodehefte nr. 2. (IO 73/6).
- Østerlund Petersen, S. 1974: Arbeidskraftundersøkelsene. Estimering av endringstall. Statistisk Sentralbyrå, Metodehefte nr. 14. (IO 74/34).