

Izlases dizaina optimizācija problemātika un daži rezultāti

Mārtiņš Liberts

Latvijas Universitāte
LR Centrālā statistikas pārvalde

23.02.2012.

Saturs

1 Problemātika

2 Mākslīgā populācija

3 Simulācijas eksperimenti

4 Dizaina efekts

5 Secinājumi

Populācija

- Galīga populācija ar apjomu N

$$U = \{1, 2, \dots, N\}$$

- Populācija ir elementu kopa

$$k \in U$$

- Katram populācijas elementam ir piekārtotas vairākas pazīmes

$$y_k, z_k, \dots$$

Populācija

- Populācijas parametri

$$Y = \sum_U y_k$$

$$\bar{Y} = \frac{1}{N} \sum_U y_k$$

$$R = \frac{Y}{Z} = \frac{\bar{Y}}{\bar{Z}}$$

Izlase

- Izlase s ar apjomu n tiek atlasīta no populācijas U ar zināmu varbūtību $p(s)$
- Izlasē iekļaušanas varbūtības katram elementam $k \in U$ ir

$$\pi_k = \sum_{s \ni k} p(s) = Pr(k \in s)$$

- Populācijas parametram var konstruēt novērtētāju $\hat{\theta}$, izmantojot s

$$\hat{\theta} = \hat{\theta}(s)$$

Izmaksas

- Apsekojuma organizācijas izmaksas

$$c_T = c_F + c_V$$

- Pieņēmums: mainīgās izmaksas ar atkarīgas no s

$$c_V = c_V(s)$$

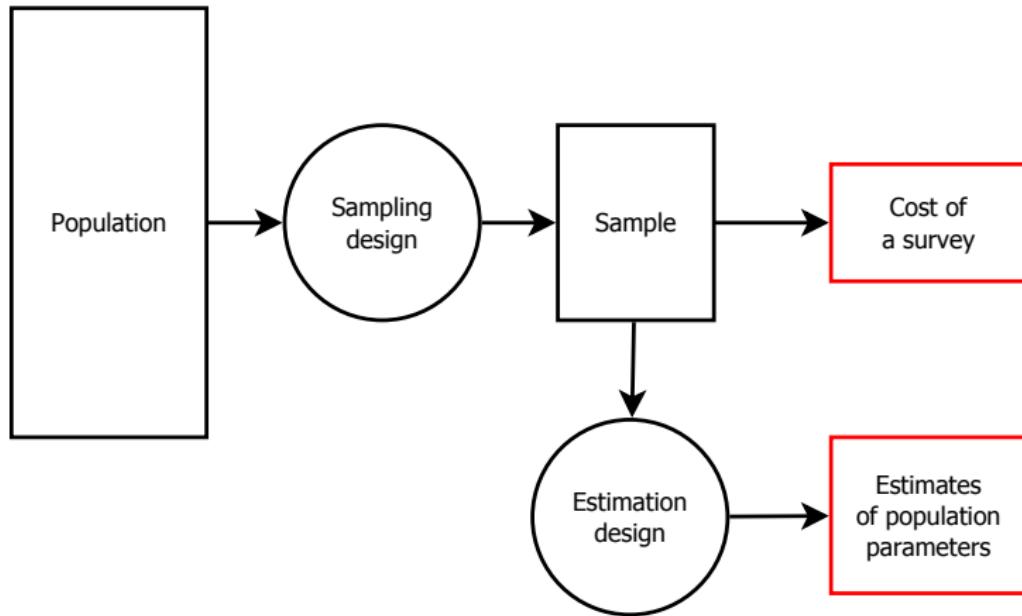
- Piemēram:

- ▶ Mainīgās izmaksas ir atkarīgas no $|s| = n$
- ▶ Mainīgās izmaksas ir atkarīgas no izlases vienību ģeogrāfiskā izvietojuma, ja apsekojums tiek veikts, izmantojot personīgās intervijas (PAPI vai CAPI)

Dizaina optimizācija

- $\hat{\theta}$ un c_V ir gadījuma lielumi, jo izlase s tiek atlasīta kā varbūtiskā izlase
- Varam pētīt $E(\hat{\theta})$, $V(\hat{\theta})$, $E(c_V)$, $V(c_V)$ īpašības
- Praksē ir vēlams sasniegt vairākus mērķus:
 - ▶ $E(\hat{\theta}) \rightarrow \theta$
 - ▶ $V(\hat{\theta}) \rightarrow \min$
 - ▶ $E(c_V) \rightarrow \min$
 - ▶ $V(c_V) \rightarrow \min$

Apsekojuma dizains



Problēma

- Praksē lietotie izlases dizaini parasti ir sarežģīti
- Tieki lietoti sarežģīti populācijas parametru novērtētāji

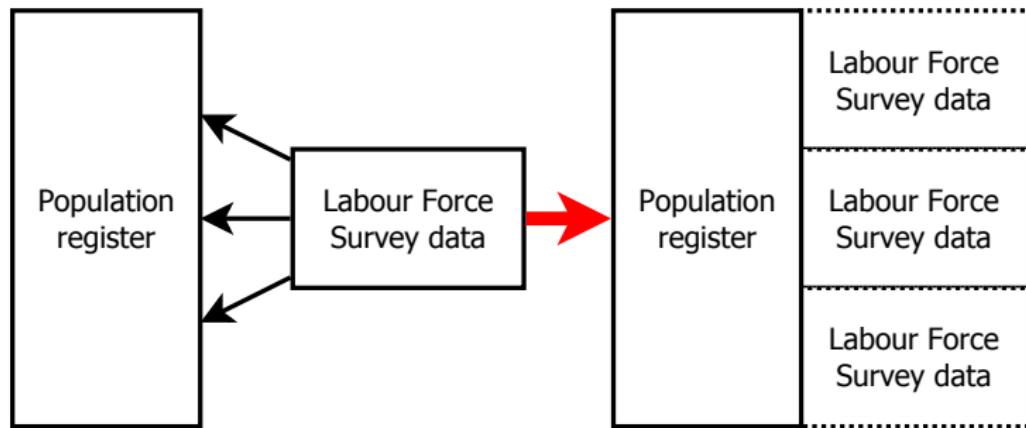
Risinājums

- iespējams risinājums – apsekojuma dizainu (izlases dizains + novērtēšanas dizains) pētīt ar simulācijas eksperimentu palīdzību
- Ir nepieciešami mākslīgas populācijas dati, lai veiktu simulācijas eksperimentus
- Mākslīgā populācija ir nepieciešama tāda, kas būtu pēc iespējas līdzīga mērķa populācijai. Līdzīga attiecībā pret:
 - ▶ informāciju, kas tiek lietota izlases veidošanai un parametru novērtēšanu;
 - ▶ ģeogrāfisko izvietojumu;
 - ▶ Pētāmajiem rādītājiem.

Mākslīgā populācija

- Mākslīgā populācija tika izveidota, izmantojot:
 - ▶ Latvijas iedzīvotāju reģistra datus;
 - ▶ Latvijas Darbaspēka apsekojuma datus.

Mākslīgā populācija



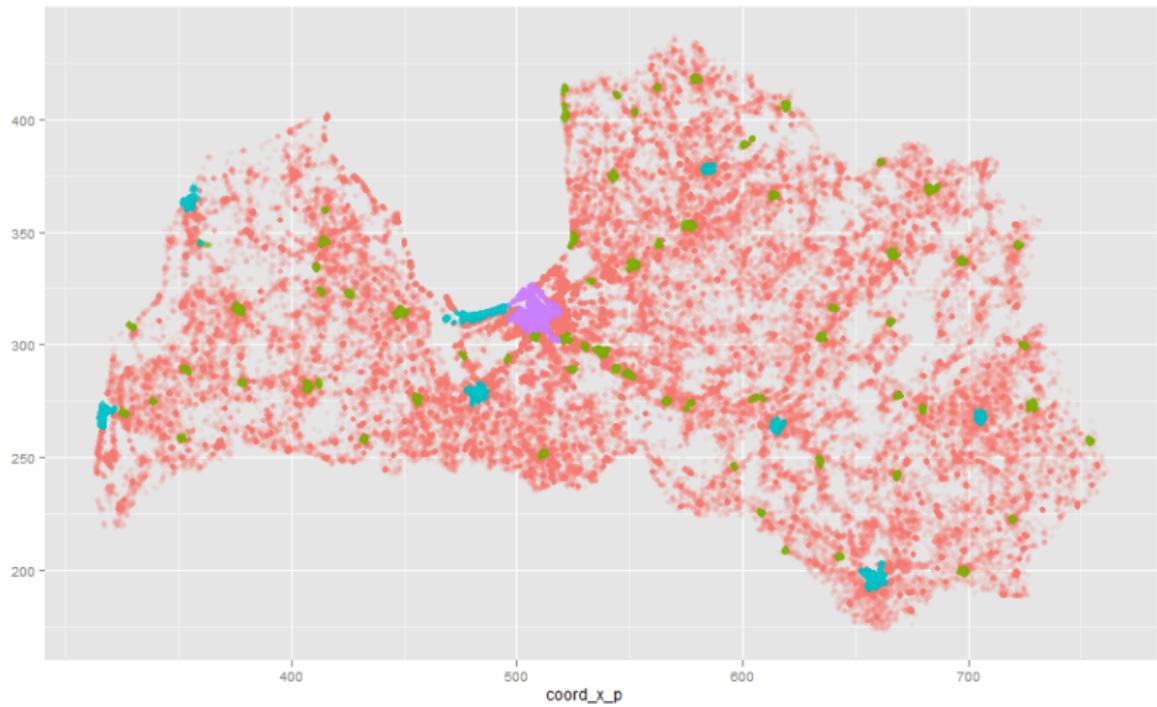
Mākslīgā populācija

Mākslīgās populācijas raksturojums

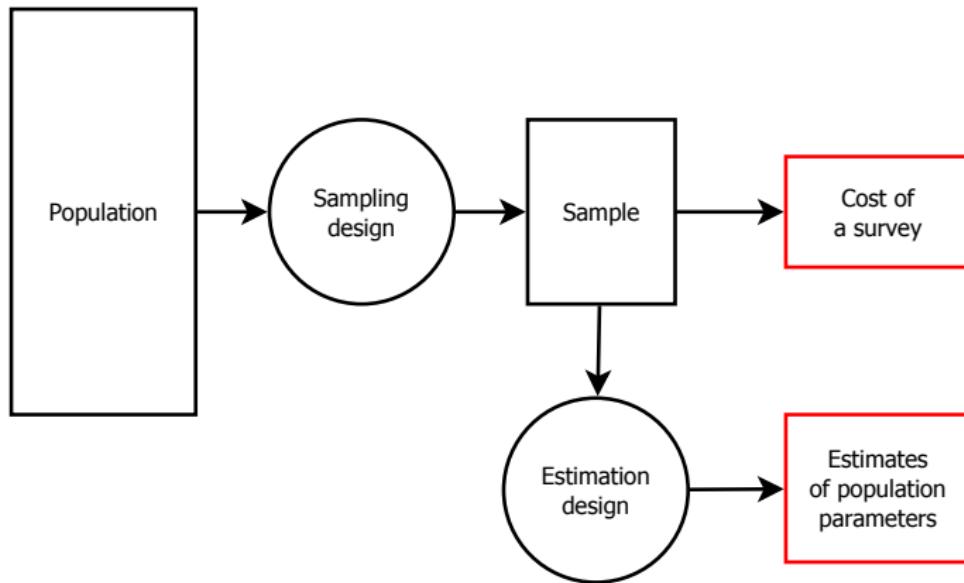
- Populācija sastāv no apmēram 1 700 000 personām (765 000 mājsaimniecībām)
- Geogrāfiskā novietojuma informācija
- Demogrāfiskā informācija (dzimums, vecums)
- Darbaspēka apsekojuma informācija par ekonomisko aktivitāti un nostrādātajām stundām

Populācijas karte

coord_y_p



Apsekojuma dizains

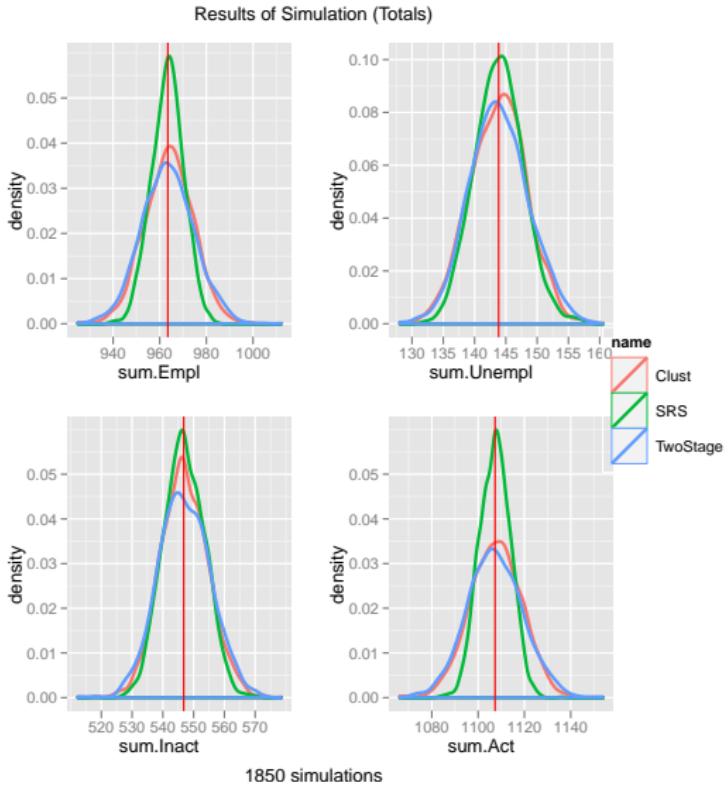


Daži rezultāti

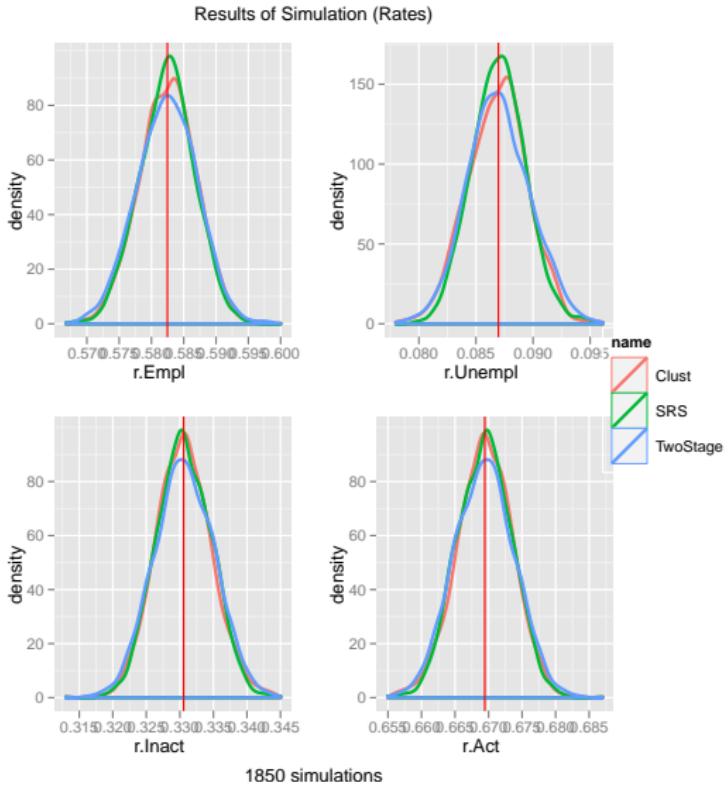
Simulācijas konfigurācija:

- Simulācijā tika izmantoti trīs izlases dizaini:
 - ▶ vienkāršā gadījuma izlase no personām ar apjomu 13 413 personas (vidēji vienā mājsaimniecībā ir 2.2 personas);
 - ▶ personu klāsterizlase, kur mājsaimniecības veido personu klāsterus, ar apjomu 6 032 klāsteri;
 - ▶ divpakāpju izlase, kur PIV ir teritorijas un SIV ir mājsaimniecības (personu klāsteri), ar apjomu 6 032 klāsteri
- HT-novērtētājs
- Simulācijas iterāciju skaits katram izlases dizainam ir 1 850
- Katrā simulācijā tika aprēķināti:
 - ▶ populācijas parametru novērtējumi;
 - ▶ ceļa garums intervētājiem, kas nepieciešams, lai veiktu apsekojumu

Rezultāti - precizitāte



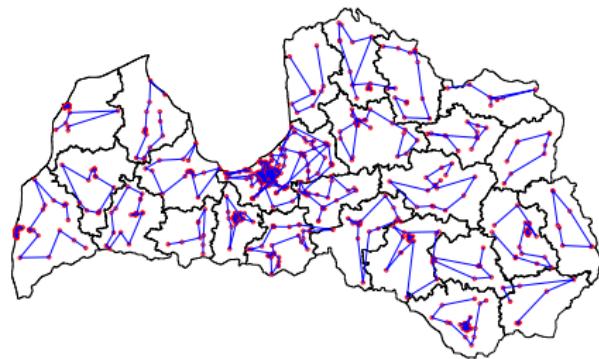
Rezultāti - precizitāte



Rezultāti - izmaksas

SRS Individuals

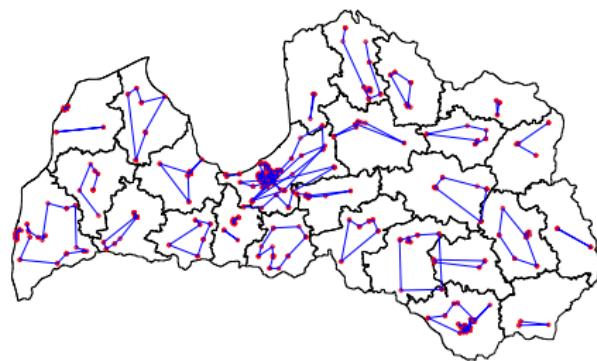
Sample size = 1032; Trip = 5048 (km)



Rezultāti - izmaksas

Cluster Sampling of persons

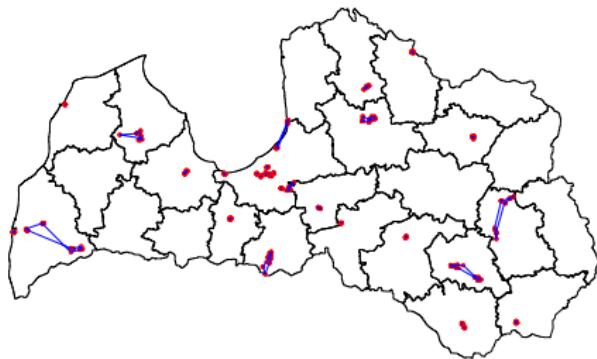
Sample size = 464; Trip = 3242 (km)



Rezultāti - izmaksas

Two Stage Sampling of Dwellings

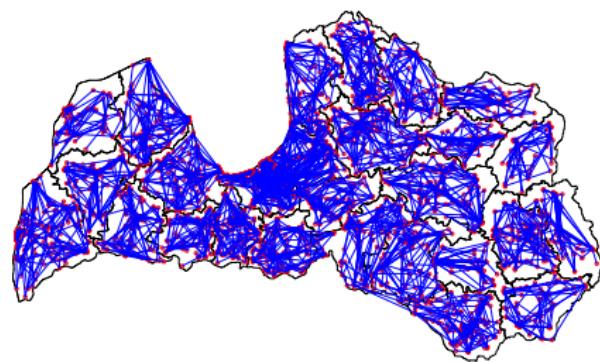
Sample size = 464; Trip = 487 (km)



Rezultāti - izmaksas

SRS Individuals

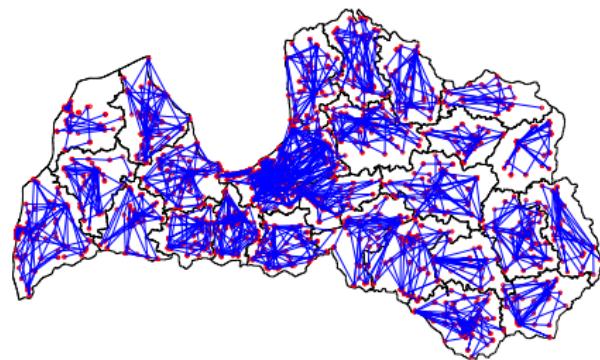
Sample size = 13413; Trip = 65312 (km)



Rezultāti - izmaksas

Cluster Sampling of persons

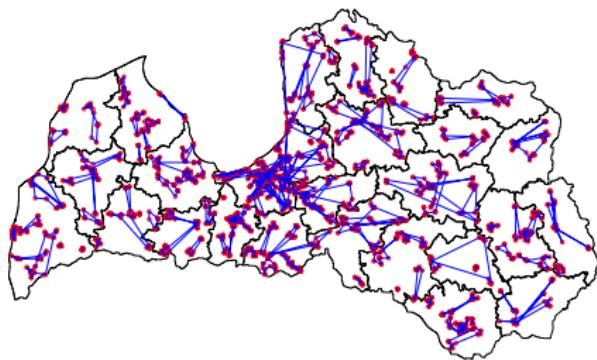
Sample size = 6032; Trip = 41357 (km)



Rezultāti - izmaksas

Two Stage Sampling of Dwellings

Sample size = 6032; Trip = 8241 (km)



Dizaina efekts

- Tradicionālais dizaina efekts – dizaina efektivitātes mērs pēc precīzitātes

$$deff_1 = \frac{Var_{Curr}}{Var_{SRS}}$$

- Varam ieviest alternatīvu dizaina efektu – dizaina efektivitātes mērs pēc izmaksām

$$deff_2 = \frac{Cost_{Curr}}{Cost_{SRS}}$$

Dizaina efekts

- Dizaina efekts – dizaina efektivitātes gan pēc precizitātes, gan pēc izmaksām

$$deff_3 = \frac{Var_{Curr} \cdot Cost_{Curr}}{Var_{SRS} \cdot Cost_{SRS}}$$

Dizaina efekts

| design | variable | var | costV | costF | cost | deff1 | deff2 | deff3 |
|----------|------------|-----|-------|-------|-------|-------|-------|-------|
| SRS | sum.Act | 45 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | sum.Act | 123 | 41963 | 32109 | 74072 | 2.737 | 0.772 | 2.113 |
| | sum.Act | 144 | 8027 | 32109 | 40137 | 3.200 | 0.418 | 1.339 |
| Clust | sum.Empl | 48 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | sum.Empl | 108 | 41963 | 32109 | 74072 | 2.221 | 0.772 | 1.715 |
| | sum.Empl | 128 | 8027 | 32109 | 40137 | 2.651 | 0.418 | 1.109 |
| TwoStage | sum.Inact | 45 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | sum.Inact | 60 | 41963 | 32109 | 74072 | 1.337 | 0.772 | 1.032 |
| | sum.Inact | 71 | 8027 | 32109 | 40137 | 1.573 | 0.418 | 0.658 |
| SRS | sum.Unempl | 15 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | sum.Unempl | 20 | 41963 | 32109 | 74072 | 1.330 | 0.772 | 1.027 |
| | sum.Unempl | 22 | 8027 | 32109 | 40137 | 1.460 | 0.418 | 0.611 |
| Clust | sum.Empl | 48 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | sum.Empl | 108 | 41963 | 32109 | 74072 | 2.221 | 0.772 | 1.715 |
| | sum.Empl | 128 | 8027 | 32109 | 40137 | 2.651 | 0.418 | 1.109 |
| TwoStage | sum.Inact | 45 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | sum.Inact | 60 | 41963 | 32109 | 74072 | 1.337 | 0.772 | 1.032 |
| | sum.Inact | 71 | 8027 | 32109 | 40137 | 1.573 | 0.418 | 0.658 |

Dizaina efekts

| design | variable | var | costV | costF | cost | deff1 | deff2 | deff3 |
|--------------------------|-------------|----------|-------|-------|-------|-------|-------|-------|
| SRS Clust TwoStage | a.Work.time | 0.012294 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | a.Work.time | 0.015167 | 41963 | 32109 | 74072 | 1.234 | 0.772 | 0.953 |
| | a.Work.time | 0.016603 | 8027 | 32109 | 40137 | 1.350 | 0.418 | 0.565 |
| SRS Clust TwoStage | r.Act | 0.000016 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | r.Act | 0.000017 | 41963 | 32109 | 74072 | 1.059 | 0.772 | 0.818 |
| | r.Act | 0.000020 | 8027 | 32109 | 40137 | 1.205 | 0.418 | 0.504 |
| SRS Clust TwoStage | r.Empl | 0.000018 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | r.Empl | 0.000019 | 41963 | 32109 | 74072 | 1.065 | 0.772 | 0.823 |
| | r.Empl | 0.000022 | 8027 | 32109 | 40137 | 1.229 | 0.418 | 0.514 |
| SRS Clust TwoStage | r.Unempl | 0.000006 | 63823 | 32109 | 95932 | 1.000 | 1.000 | 1.000 |
| | r.Unempl | 0.000007 | 41963 | 32109 | 74072 | 1.233 | 0.772 | 0.952 |
| | r.Unempl | 0.000008 | 8027 | 32109 | 40137 | 1.370 | 0.418 | 0.573 |

Dizaina efekts

- Tradicionālais dizaina efekts – nosacījums: vienāds izlases apjoms

$$deff_1 = \frac{Var_{Curr}(n_{Curr} = n)}{Var_{SRS}(n_{SRS} = n)}$$

- Alternatīvs dizaina efektu – nosacījums: vienādas izmaksas

$$deff_4 = \frac{Var_{Curr}(c_{Curr} = c)}{Var_{SRS}(c_{SRS} = c)}$$

Secinājumi

- Simulācijas eksperimenti – rīks, lai izvērtētu un analizētu apsekojuma dizaina īpašības:
 - ▶ precīzitāte
 - ▶ izmaksas
- Noderīgs:
 - ▶ analizējot apsekojuma efektivitāti
 - ▶ plānojot izmaiņas apsekojuma dizainā
 - ▶ veicot izmaiņas apsekojuma dizainā

Paldies par uzmanību!