

Natalja Lepik ja Imbi Traat (Tartu Ülikool), 2013



E-kursuse „Valikuuringute teooria I” materjalid

Aine maht 6 EAP

Natalja Lepik ja Imbi Traat (Tartu Ülikool), 2013

Sisukord

1	Sissejuhatus, mõisted	4
1.1	VU aine põhijooned	5
1.2	Valikumeetodid	6
1.3	Mõisteid valikuuringute praktikast	8
1.4	Erinevus mat. statistikast. VU teooria eripära	9
2	ÜK parameetrid. Valikudisain, selle karakteristikud	10
2.1	ÜK parameetrid	10
2.2	Tunnuse dispersioon ÜK-s	11
2.3	Valikudisain	11
2.3.1	TTA-disainid, (<i>WOR - Without Replacement</i>)	14
2.3.2	Näiteid lihtsa funktsionaalse kujuga TTA-disainidest	14
2.3.3	TGA-disainide näited	17
2.4	Disaini karakteristikud	20
3	Hindamise alused	22
3.1	Valikuuringu andmed	22
3.2	Kogusumma nihketa hindamine	24
3.2.1	Dispersiooni hindamine fikseeritud mahuga disainide korral	28
3.2.2	Nihketa hinnang kogusummale TGA disainide korral	30
3.2.3	Üldkogumi keskmise nihketa hindamine	32
3.2.4	Kahe nihketa hinnangu kovariatsioon	36
3.3	Suhte hinnang	37
3.3.1	Taylori rida kahe argumendi korral	37
3.3.2	Suhte hinnangu Taylori rittaarendus	37
3.4	Osakogumi hindamine.	39
3.5	Hinnangu täpsuse iseloomustamine	43
3.5.1	Valimimahu määramine	44
3.6	Disainiefekt	45
4	Hindamine lihtsa juhuvaliku korral, TTA	45
4.1	LJV disainikarakteristikud	46
4.2	Hinnang kogusummale LJV korral	47
4.3	Kovariatsioon kahe hinnangu vahel LJV TTA korral	51
4.4	Suhtehinnang LJV TTA korral	52

<i>Valikuuringute teooria I, Natalja Lepik, Imbi Traat; 2013</i>	3
4.5 Hindamine osakogumites LJV TTA korral	53
5 Hindamine lihtsa juhuvaliku TGA korral	56
5.1 Isekaaluvad disainid	60
6 Süstemaatiline valik	61
6.1 Hindamine SÜ korral	63
6.2 SÜ disaini efekt	65
6.3 SÜ realiseerimine praktikas	66
7 Ebavõrdsete tõenäosustega valik	67
7.1 Suurusega võrdelise tõenäosusega valik	69
7.2 Poissoni valik	70
8 Kihtvalik	71
8.1 Hindamine kihtvaliku korral	72
8.2 Lihtne juhuslik kihtvalik	74
8.3 Valimi optimaalne paigutus	77
8.4 Optimaalne valimi paigutus KLJV korral	80
8.5 Alternatiivsed valimi paigutused KLJV korral	81
8.6 LJV ja KLJV võrdlemine	82
9 Järeldamine	84
9.1 Järeldamine LJV korral	86
10 Klastervalik	86
11 Kahe-astmeline valik	87
11.1 Tähistused	87
11.2 Hindamine kahe-astmelise valiku korral	88
11.3 Kahe-astmeline lihtne juhuslik valik	91
11.4 Isekaaluv kahe-astmeline valik	91
12 Abiinformatsiooni kasutamine hinnangutes	93
12.1 Regressioonimudel üldkogumi jaoks	94
12.2 Regressioonihinnang	95

Hindamise kriteeriumid

Hindamine

Aine lõpeb eksamiga, eksamile pääsemiseks peab olema arvestatud projekt ja loengutel antud ülesanded.

Aine hinne koosneb: 0,4 referaat + 0,6 eksam,

kusjuures eksam peab olema sooritatud vähemalt 51% ulatuses.

Märkus loengutel antavate ülesannete kohta

Igal loengul antakse 3 koduülesannet, millest ainult ühe peab esitama järgmiseks loenguks (omal vabal valikul).

Juhul, kui tähtaeg saab ületatud, võib seda pikendada ühe nädala võrra lahendades juba kaks ülesannet. Kui teine tähtaeg saab ka ületatud, siis ootan päev enne eksamit kõik kolm ülesannet ja seda iga loengu kohta.

Kirjandus

1. W. G. Cochran (1977) Sampling Techniques. Third edition. Wiley
2. C.-E. Särndal, B. Swensson, J. Wretman (1992) Model Assisted Survey Sampling, Springer-Verlag
3. I. Traat, J. Inno (1997) Tõenäosuslik valikuuring, TÜ Kirjastus
4. Loengute konspektid **moodles**.

1 Sissejuhatus, mõisted

Üldkogum e. **populatsioon** (*population*) - objektide hulk (lõplik hulk), mille kohta soovitakse vastavalt püstitatud probleemülesandele saada informatsiooni.

Osakogum (*subpopulation*) - üldkogumi alamhulk, mis on fikseeritud tausttunnuse või uuritava tunnuse väärtuste järgi ja mida soovitakse eraldi uurida. Osakogumi objektid on sama tüüpi, mis üldkogumi omad (pered mõlemas, isikud mõlemas vms.)

Andmekogumismeetodid:

- **Kõikne uuring** e. **loendus** - andmete kogumine üldkogumi kõigilt objektidelt.
 - "+": võimaldab täpset infot ÜK kohta fikseeritud ajahetkel.
 - ": töömahukas, kallis.
 - ": mahukuse tõttu kumuleeruvad vead.

- **Register** - andmebaasid mitmesuguste ÜK-te kohta, nt. rahvastikuregister, äriregister, hooneregister jne.
"+": regulaarselt täpsustatud andmed aruannete põhjal.
"-": registritest saadavad tunnused on fikseeritud registri ülesehitusega.
- **Valikuuring** - statistiline uuring, milles otsustused üldkogumi kohta tehakse valimi baasil.

VU eelised:

- väiksem maksumus
- suurem kiirus
- paindlikkus
- laiem rakendatavus
- suurem täpsus

Valikuuringute teooria - teadus andmete kogumisest, töötlemisest, esitamisest ja analüüsimisest. Tegevus algab probleemülesande püstitamisest ja lõpeb tulemuste publitseerimisega.

Esimene VU kursus Eestis - 1993.

1.1 VU aine põhijooned

- Valikustrateegiad:
 - valimi võtmise meetod (lõpmata palju viise, põhjustavad erinevuse hinnangutes sama parameetri kohta)
 - hinnangufunktsiooni valik

Eesmärgid:

- Leida selline strateegia, mis minimiseeriks hinnangu dispersiooni (juh. viga) ja/või minimiseeriks uuringu maksumuse.
- Leida praktikas realiseeritav strateegia, osata anda nihketa hinnanguid ja nende täpsushinnanguid.

- Muude vea allikate analüüs (näiteks kadu), vastavad kompenseerimis-meetodid.
- Lisainformatsiooni kasutamine (registrite info), et muuta hinnanguid täpsemaks, ka kooskõlaliseks muude teadaolevate näitajatega.

1.2 Valikumeetodid

Tõenäosuslikud - iga ÜK objekti jaoks on teada tema valimisse sattumise (kaasamise) tõenäosus. Läheb vaja ka 2 objekti koos valimisse sattumise tõenäosust, mis aga iga kord pole täpselt leitav.

Tõenäosuslikud meetodid jagunevad laias laastus:

- tagasipanekuta valik – TTA (*sampling without replacement*); ÜK iga objekt saab olla valitud max 1 kord
- tagasipanekuga valik – TGA (*sampling with replacement*); ÜK iga objekt saab olla valimis rohkem kui 1 kord

Objekti valimise viisi järgi jagatakse meetodeid ka tõmbe- või loeteluviisi valikuteks.

Levinuimad tõenäosuslikud valikumeetodid:

- **Lihtne juhuslik valik.** On olemas TTA ja TGA. Kõikidel objektidel on võrdne valimisse sattumise tõenäosus. Veelgi enam, kõik valimid on võrdse esinemistõenäosusega. Tulemuseks on valim etteantud mahuga n .
- **Poissoni valik.** Iga ÜK objekti valimisse kaasamine otsustatakse sõltumatult Bernoulli juhusliku suuruse abil, kusjuures igal objektil võib olla erinev kaasamistõenäosus, π , $0 < \pi < 1$. Tulemusena saadud valim on juhusliku valimimahuga, $E\mathbf{n} = \sum_U \pi_i$. Juh. valimimaht suurendab hinnangute varieeruvust. Võrdsete kaasamistõenäosuste erijuhul on tegemist Bernoulli valikuga.
- **Süsteemiline valik** (*Systematic sampling*). Objektide valimine toimub fiks. sammu tagant loendist, kusjuures esimene valitav objekt määratakse juhuslikult. Tulemuse täpsus sõltub objektide paigutusest loendis.

- **Suurusega võrdelise tõenäosusega valik** (*Sampling with probabilities proportional to size*). Objektide kaastamistõenäosused on võrdelised objektide suurustega. Saavutatakse suurte objektide ülesindatus valimis. Hinnangute arvutamisel tuleb andmeid erinevalt kaaluda. Tulemusena tõuseb mõnede hinnangute täpsus.
- **Kihtvalik** (*Stratified sampling*). ÜK jagatakse osadeks ehk kihtideks. Igas kihis rakendatakse sõltumatult mingit tõen. valikumeetodit. Igas kihis arvutatakse parameetrite hinnanguid, neid sobivalt kombineerides saab leida ÜK hinnanguid. Kihid peavad olema uuritava tunnuse suhtes võimalikult homogeensed.
- **Klastervalik** (*Cluster sampling*). ÜK koosneb objektgruppidest ehk klastritest. Toimub klastrite juhuslik valik. Parameetrid arvutatakse igas klastris kõiki objekte kasutades.
- **Kaheastmeline valik** (*Two-stage sampling*). Esimesel sammul toimub klastervalik ja teisel – igas klastris toimub objektide juhuslik valik. Klastrid peavad olema võimalikult heterogeensed.

Empiirilised valikud - kaasamistõenäosusi pole teada; eesmärgiks on saada ÜK struktuuriga sarnane valim. Pole võimalik leida kvantitatiivseid täpsusnäitajaid.

Empiirilise valiku meetodid:

- **Kvootide meetod** (*quota sampling*). Valimi struktuur määratakse tausttunnuste järgi (sagedus)
- **Expertvalik** (*Expert sampling*). Subjektiivse valiku teostab ekspert.
- **Tasakaalustatud valik**. Valik on sarnane kvootide valikule, kuid valimi struktuuri määravad mite tausttunnuste osakaalud, vaid muud näitajad, nt keskmine (vanus).

Märkus. Suhteliselt hiljuti on välja töötatud tõenäosuslik tasakaalustatud valik (De Ville). Tulemusena saadakse juhuslik valim, milles objektide kaastamistõenäosused on teada.

1.3 Mõisteid valikuuringute praktikast

Objekt, element, ühik, indiviid.

Tunnus - uuritav või taust- ehk abitunnus.

Üldkogum, populatsioon.

Osakogum - ÜK alamhulk, mis on fikseeritav tausttunnuse või uuritava tunnuse väärtuste järgi.

Loend, freim (*frame*) - ÜK elementide loend, mis koosneb ÜK elementidest või nende gruppidest.

Freimi abil peab olema võimalik...

- (1)... valida valimit vastavalt fikseeritud valikudisainile
- (2)... saada kontakti valitud ÜK elementidega

Eristatakse **kahte liiki ÜK-meid**:

1. **sihtkogum** (*target population*) - objektide hulk, mis tuleb uurida lähitavalts statistilisest ülesandest
2. **loendile vastav ÜK** (*frame population*)

Aktuaalne kogum - objektide hulk, mis kuulub nii loendisse kui sihtkogumisse.

Valim - aktuaalse kogumi osahulk, mis määratakse statistilise valikumeetodiga.

Loendi vead:

- **ülekaetus** - sisaldab ka ÜK-sse mitte-kuuluvaid elemente
- **alakaetus** - ei sisalda kõiki ÜK-sse kuuluvaid elemente
- **kordumised** - mõni ÜK-i element on kirjeldatud mitmel korral

Näide: Ettevõtete uuring

ülekaetus - loendis on tegevuse lõpetanud ettevõtted

alakaetus - äsja tegevust alustanud ettevõtted

Isekaaluv valim - ühesuguse tähtsusega objektidest koosnev valim, iga objekt valimis esindab võrdse arvu ÜK objekte. Hinnangute arvutamisel isekaaluvalt valimilt ei ole objektidele vaja omistada erinevaid kaalusid.

Kadu - valimi osa, mis mingil põhjusel jääb uuringust kõrvale.

Kao määr - kao osakaal valimist.

Vastamismäär - vastanute osakaal valimist.

1.4 Erinevus mat. statistikast. VU teooria eripära

Klassikaline mat. statistika:

1. ÜK on lõpmatu. Kui ta ongi lõplik, siis valik on tagasipanekuga, mistõttu valimimaht võib ikka lõpmatult kasvada.
2. Juh. suuruse Y käitumine on kirjeldatud jaotusega.
3. Juh. suurus koos oma jaotusega annab ÜK mudeli: $Y \sim F(\theta)$. Tahame hinnata parameetrit θ .
4. Juh. valimi element y_i on juh. suuruse Y realisatsioon. Realisatsioonid pärinevad sõltumatult samast jaotusest (ssj)
5. ssj-eeldus lubab leida parameetrile hinnanguid $\hat{\theta} = \hat{\theta}(y_1, y_2, \dots, y_n)$ ja uurida $\hat{\theta}$ statistilisi omadusi.

Valikuuringute teooria:

1. ÜK on reaalne, lõplik: $U = 1, 2, \dots, N$.
2. TTA valiku korral ei saa valimimaht lõpmatult suurened.
3. Tunnuse väärtusi y_1, y_2, \dots, y_N võib küll vaadelda diskreetse üldkogumi-jaotusena, kus $p(y_i) = \frac{1}{N}, \forall i$, kuid TTA valiku korral ei toimu valik igal sammul samast ÜK jaotusest.
4. Realiseerunud väärtused y_i pole ssj (va LJV TGA valiku korral).
5. Üldjuhul valim ei peegelda ÜK-t, st hindamisprobleemid vajavad teist lähenemist
6. Hinnangute omadused on määratud valikudisainiga – valimite tõenäosustega.

⇒ Teatud mõttes on klassikaline statistika vaadeldav VU osana! LJV TGA valim on nagu klassikalise statistika valim, kus kehtivad klassikalise statistika tulemused.

2 ÜK parameetrid. Valikudisain, selle karakteristikud

2.1 ÜK parameetrid

Olgu lõplik ÜK $U = \{1, 2, \dots, N\}$, N - ÜK maht ja uuritav tunnus y .

Def. ÜK parameetriks θ nim. arvulist näitajat, mis mingis mõttes iseloomustab üldkogumit.

Näiteid:

- uuritava tunnuse **kogusumma** (*total*): $t = \sum_{i=1}^N y_i = \sum_U y_i$. Kui on vaja eristada kahe tunnuse kogusummasid, siis lisame indekseid, t_y või t_x . Mõnes kirjanduses kasutatakse vastavalt tähiseid Y ja X ;
- ÜK **keskmine** (*mean*): $\bar{Y} = t_y/N = t_y / \sum_U 1$;
- osakogumi (*domain*) $U_d \subset U$ **osakaal**: $P_d = N_d / N$, kus N_d on osakogumi U_d maht;
- kahe kogusumma **suhe** (*ratio*): $R = t_y / t_x = \sum_U y_i / \sum_U x_i$.

Kõik eespool nimetatud parameetrid avalduvad kogusummade kaudu!

Osakogumi U_d jaoks defineerime binaarse tunnuse z , kus

$$z_i = \begin{cases} 1 & , \text{ kui } i \in U_d; \\ 0 & , \text{ muidu.} \end{cases}$$

Osakogumi U_d maht N_d on tunnuse z kogusumma:

$$N_d = t_z = \sum_U z_i = \sum_{U_d} 1,$$

ja U_d osakaal P_d ÜK-s on tunnuse z keskmine:

$$P_d = \bar{Z} = \frac{t_z}{N}.$$

Järeldus. VU teooria pühendab oma tähelepanu kogusumma hindamisprobleemile.

2.2 Tunnuse dispersioon ÜK-s

Tunnuse y dispersiooniks nim. suurust:

$$S_{y,U}^2 = \frac{1}{N-1} \sum_U (y_i - \bar{Y})^2 = \frac{1}{N-1} \left[\sum_U y_i^2 - N\bar{Y}^2 \right] = \frac{1}{N-1} \left[t_y^2 - \frac{t_y^2}{N} \right].$$

Binaarse tunnuse z korral saab selle valemi teisendada järgmisele kujule:

$$S_{z,U}^2 = \frac{N}{N-1} P_d Q_d,$$

kus $Q_d = 1 - P_d$.

Näidata!

2.3 Valikudisain

Olgu $U = 1, 2, \dots, N$.

Valimi esitamise viisid:

1) **Hulkvalim**, s - ÜK-i osahulk, $s \in U$, elementide järjestus pole tähtis. Kokku on võimalik saada 2^N hulka (*Miks?*).

Kui palju on võimalik saada valimeid mahuga n ?

Näide: $s = \{2, 5, 3\}$ - korduvaid elemente pole!

Kasutatakse TTA disainide koral; kõige rohkem levinud.

2) **Järjestusvalim**, $js = \{i_1, i_2, \dots, i_n\}$, $i_k \in U$ – valimi elemendid on esitatud elemendi võtmise järjekorras; võivad esineda ka kordused.

Näiteks: $js = \{3, 1, 5, 1\}$.

Kasutatakse TGA disainide korral; praktikas pole levinud, kuid teooria on nende jaoks lihtsam.

3) **Vektorvalim**, $k = (k_1, k_2, \dots, k_N)$, sama dimensiooniga nagu ÜK, kus k_i on objekti i valikute arv. Juhul, kui $k_i = 0$, siis objekt i pole valimis.

Näiteks, $k = (1, 0, 2, 0, \dots, 3)$.

Valimimaht: $n = \sum_{i=1}^N k_i$.

Saab kasutada nii TTA kui ka TGA disainide korral. Kasutame selles kursuses teooria arendustes.

Def. Juhuslikku vektorit $I = (I_1, I_2, \dots, I_N)$ nim. valikuvektoriks, kus I_i (valikuindikaator) näitab objekti i valikute arvu ($i \in U$).

Valikuvektori realisatsiooniks on (vektor)valim k .

NB! $\mathbf{n} = \sum_{i=1}^N I_i = \sum_U I_i$ - valimimaht, mis võib olla juhuslik.

Def. Valikudisainiks nim. valikuvektori I jaotust:

$$I \sim p(k), p(k) = P(I = k), \sum_k p(k) = 1.$$

Toodud definitsioon on matemaatiliseks aluseks teooria arendamisel. Meie tegeleme selles kursuses nn disainipõhise lähenemisega, kus hinnangute omadused on määratud valikudisainiga.

Praktikas räägitakse valikudisainist ka kui reeglite kogumist, kuidas valimit üldkogumist võtta. Lõppkokkuvõttes määrab aga ka reeglite kogum üheselt võimalike valimite tõenäosused ehk valikudisaini.

Def. Valikudisaini nim. fikseeritud mahuga n disainiks, kui selle disaini korral $\sum_U I_i \equiv n$.

Def. Valikudisaini nim. tagasipanekuta disainiks, TTA, kui $I_i \in \{0, 1\} \forall i$, muidu tagasipanekuga, TGA.

Def. ÜK objekti i ($i = 1, \dots, N$) kaasamistõenäosuseks π_i nim.tõenäosust, millega see objekt kaasatakse valimisse antud disaini korral:

$$\pi_i = P(i \in s) = P(I_i \geq 1) = \sum_{k, k_i \geq 1} p(k).$$

Erijuhul, kui tegemist on TTA disainiga, siis $\pi_i = P(I_i = 1)$.

Näide. TTA, $N = 4$, $n = 2$.

U	Valimid veerus						π_i
	1.	2.	3.	4.	5.	6.	
1	1	1	1	0	0	0	$\pi_1 = 0,6$
2	1	0	0	1	1	0	$\pi_2 = 0,6$
3	0	1	0	1	0	1	$\pi_3 = 0,4$
4	0	0	1	0	1	1	$\pi_4 = 0,4$
$p(k)$	0,4	0,1	0,1	0,1	0,1	0,2	$\sum_k p(k) = 1$

Valimimaht: $n = \sum_{i=1}^4 \pi_i = 2$.

Def. Objektide i, j 2. järku kaasamistõenäosuseks nim. tõenäosust, millega need objektid kaasatakse korraga valimisse antud disaini korral:

$$\pi_{ij} = P(I_i \geq 1, I_j \geq 1).$$

Leida π_{12} eelmise näite jaoks!

Valikuindikaatori omadused TTA disaini korral:

- $E(I_i) = \pi_i$;
- $V(I_i) = \pi_i(1 - \pi_i)$;
- $Cov(I_i, I_j) = \pi_{ij} - \pi_i\pi_j$.

Näidata!

2.3.1 TTA-disainid, (*WOR - Without Replacement*)

Sel juhul on I_i Bernoulli juh. suurus:

$$I_i \sim Bin(1, \pi_i) = Be(\pi_i), \quad I_i \in \{0, 1\},$$

$$\pi_i = P(I_i = 1),$$

mille jaotusfunktsioon avaldub järgmiselt:

$$P(I_i = k_i) = \pi_i^{k_i} (1 - \pi_i)^{1-k_i}, \quad k_i \in \{0, 1\}.$$

Def. Juh. vektorile $I = (I_1, I_2, \dots, I_N)$ vastavat jaotust TTA disainide korral nim. mitmemõõtmeliseks Bernoulli jaotuseks (MB).

Jaotuste keeles on kõik TTA disainid on MB erijuhtumid!

MB jaotust iseloomustab:

- pole üldist funktsionaalset vormi;
- on võimalik ette anda kõikvõimalike tõenäosuste tabelina:

$$P(I = k) = p(k), \quad k \in \{0, 1\}^N,$$

$$0 \leq p(k) \leq 1, \quad \sum p(k) = 1;$$

- $p(k)$ jaoks on võimalik lõpmata palju erinevaid variante!
- mitmed klassikalised valikudisainid on lihtsa funktsionaalse kujuga MB jaotused.

2.3.2 Näiteid lihtsa funktsionaalse kujuga TTA-disainidest

1. Poissoni disain

$$\left[\begin{array}{l} I = (I_1, \dots, I_N), \quad I_i \perp I_j, \quad i \neq j; \\ I_i \sim Be(\pi_i); \\ I \sim p(k) = P(I = k) = \prod_{i=1}^N P(I_i = k_i) \\ \quad = \prod_{i=1}^N \pi_i^{k_i} (1 - \pi_i)^{1-k_i}. \end{array} \right.$$

Poissoni valimi gener. algoritm: (gener. N korda sõltumatult Bernoulli juh. suuruseid)

$$i = 1;$$

$$\left[\begin{array}{l} u \sim U(0, 1); \\ \text{if } u < \pi_i \text{ then } I_i = 1 \text{ else } I_i = 0; \\ i := i + 1. \end{array} \right.$$

2. Bernoulli disain

$$\left[\begin{array}{l} \text{Poissoni disaini erijuht, kus } \pi_i \equiv \pi; \\ p(k) = \pi^{|k|}(1 - \pi)^{N - |k|}, \text{ kus } |k| = \sum_{i=1}^N k_i. \end{array} \right.$$

Näidata, et $\mathbf{n} \sim \text{Bin}(N, \pi)$.

3. Lihtne juhuslik valik TTA

$$\left[\begin{array}{l} I \sim p(k), I_i \not\sim I_j; \\ p(k) = \begin{cases} \frac{1}{C_N^n}, & \text{kui } |k| = n, C_N^n = \frac{N!}{(N-n)!n!}; \\ 0, & \text{muidu.} \end{cases} \end{array} \right.$$

Kõikidel valimitel mahuga n on võrdne tõenäosus olla valitud.

Valimi genereerimise võimalused (palju):

(i) Definitsiooni järgi. Loetleda kõikvõimalikud valimid mahuga n (selliseid võimalusi on C_N^n) ja siis valida üks valim võrdse tõenäosusega, nt. urnist.

(ii) Tõmbeviis (elemendid on nummerdatud)

$$\left[\begin{array}{l} i = 1 \Rightarrow \text{valime tn-ga } 1/N \\ \text{ja eemaldame } \ddot{U}K\text{-st;} \\ \\ i = 2, \dots, n \Rightarrow \text{valime tn-ga } 1/[N - (i - 1)] \\ \text{ja eemaldame iga kord } \ddot{U}K\text{-st.} \end{array} \right.$$

Arvutis saab valitava elemendi kätte järgmise eeskirja abil:

$$i.\text{nda el. nr} = \lfloor (N - i + 1) \cdot U(0, 1) + 1 \rfloor,$$

kus $\lfloor a \rfloor$ tähistab arvu a täisosa.

(iii) Loeteluviis (tulemuseks on vektorvalim)

$\forall i = 1, \dots, N$ seame vastavusse juh. arvu $u_i \sim U(0, 1)$.

$$\left[\begin{array}{l} i = 1 : \text{ kui } u_1 < n/N \Rightarrow \text{1. el. on valimis} \\ \\ i = 2, \dots, N : \\ \quad \text{kui } u_i < \frac{n-n_i}{N-i+1} \Rightarrow i.\text{s el. on valimis.} \end{array} \right.$$

Siin n_i on el-ntide arv, mis on valitud ÜK-i esimese $i - 1$ objekti seast.

(iv) Järjestusvalik

$$\left[\begin{array}{l} \forall i = 1, \dots, N \text{ seame vastavusse juh. arvud } u_1, \dots, u_N, u_i \sim U(0, 1). \\ \text{Järjestame ÜK objektid ümber } u_i \text{ järgi kasvavalt: } u_{(i_1)} < u_{(i_2)} < \dots < u_{(i_N)} \\ \text{Võtame valimisse esimest } n \text{ objekti.} \end{array} \right.$$

NB! Saab ka kasutada suvalist pidevat jaotust!

Tegelikult on iga n elemeline komplekt, mis on võetud järjestatud failist LJV TTA valim. Seda omadust saab valikuuringutes kasutada vastamiskoormuse reguleerimiseks.

4. Tinglik Poissoni disain

Olgu $I \sim p(k) = \prod_{i=1}^N \pi_i^{k_i} (1 - \pi_i)^{1-k_i}$ Poissoni disain.

Tinglik Poissoni disain:

$$I^{TP} \sim P(I = k | \sum_U I_i = n) = \begin{cases} \frac{p(k)}{P(\sum_U I_i = n)}, & \text{kui } \sum_{i=1}^N k_i = n; \\ 0, & \text{muidu.} \end{cases}$$

Tingliku Poissoni valimi genereerimine:

$$\left[\begin{array}{l} \text{Teostada Poissioni valik nii nagu on kirjeldatud 1. näites.} \\ \text{Kui valimimaht pole } n, \text{ siis jätta saadud valim kõrvale ja alustada uuesti.} \\ \text{Korrata nii kaua, kuni saavutatakse vajalik valimimaht.} \end{array} \right.$$

2.3.3 TGA-disainide näited

$$I = (I_1, \dots, I_N) \sim p(k) = P(I = k),$$

$$I_i = \begin{cases} 0, & i \text{ ei ole valimis;} \\ k_i > 0, & i \text{ on valimis } k_i \text{ korda, } k_i \in \{0, 1, 2, \dots, n\}. \end{cases}$$

1. Multinomiaaldisain:

- Valikutõenäosused p_i on fikseeritud iga i jaoks, $i \in U$ kogu valikuprotsessis, $\sum_{i=1}^N p_i = 1$.
- Objekt valitakse vastavalt p_i -le, registreeritakse ja seejärel pannakse tagasi ÜK-sse.
- Protsessi korratakse n korda (kuni valim on käes).

Tähistame: $I \sim M(n; p_1, p_2, \dots, p_N)$, mille tõenäosusfunktsioon on järgmine:

$$p(k) = \frac{n!}{\prod_{i=1}^N k_i!} \prod_{i=1}^N p_i^{k_i}, \text{ kui } |k| = n.$$

Erijuht

Juhul, kui kõik p_i on võrdsed, $p_i \equiv \frac{1}{N}$, siis on tegemist lihtsa juhuvalikuga TGA:

$$p(k) = \frac{n!}{N^n \prod_{i=1}^N k_i!}, \text{ kui } |k| = n.$$

- M() disaini korral: $I_i \sim B(n, p_i)$, st et objekt i saab olla valitud $k_i = 1, \dots, n$ korda.

Kirjuta ise välja $E(I_i)$ ja $V(I_i)$!

- LJV TGA korral: $I_i \sim B(n, \frac{1}{N})$, $E(I_i) = \dots?$

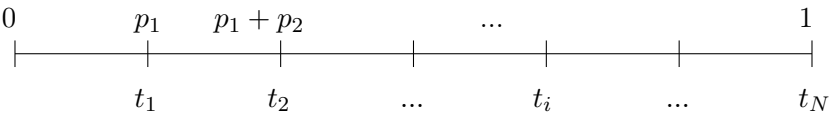
Mult.disaini genereerimine

Kasutame nn kumulatiivsete summade meetodit...

Moodustada kum. summad:

$$t_i = \sum_{j=1}^i p_j, \quad i = 1, \dots, N.$$

Need summad asuvad lõigul $[0, 1]$:



Genereerida $u \leftarrow U(0, 1)$

Kui $u \in (t_{i-1}, t_i]$ siis element i on valimis.

Korrata protseduuri n korda.

Näide R-is. LJV TGA genereerimine

```

N=20 # ÜK maht
Nr=1:N # ÜK  $U = (1, \dots, 20)$ 
n=10 # valimi maht

# Lihtne juhuvalik TGA
p=rep(1/N,N) # valikutõenäosuste vektor  $1 \times N$ 
cum=cumsum(p) # kumulatiivsete summade vektor  $1 \times N$ 
s=rep(NA,n) # valimi vektor  $1 \times n$ 

for(i in 1:n){
  u=runif(1) # juhuslik arv  $\leftarrow U(0,1)$ 
  j=1 # hakkame otsima, mis lõiku  $u$  kuulub
  while(u>cum[j]){
    j=j+1}
  s[i]=Nr[j]}
sort(s) # järjestame andmeid valimis

```

Mõned tulemused:

$s1$: 1, 7, 8, 8, 11, 13, 14, 14, 17, 20

$s2$: 6, 7, 8, 11, 13, 13, 15, 18, 18, 19

$s3$: 2, 7, 7, 8, 8, 9, 13, 14, 15, 20

2. Hüpergeomeetriline disain

Iga element saab olla valitud kuni m_i korda:

$$I_i \in \{0, \dots, m_i\}, \quad i = 1, 2, \dots, N; \quad m_i < n.$$

Olgu $m = \sum_{i=1}^N m_i$.

Tähistame $I \sim HG(n; m_1, m_2, \dots, m_N)$, kus jaotuse tõenäosusfunktsioon on järgmine:

$$p(k) = P(I = k) = \frac{\prod_{i=1}^N C_{m_i}^{k_i}}{C_m^n}, \quad \text{kui } |k| = n.$$

Erijuht. Kui kõik $m_i \equiv 1$, siis annab HG disain LJV TTA:

$$p(k) = \frac{1}{C_N^n}.$$

Harjutusülesanne (teeme loengul)

Olgu $N = 4$, $E\mathbf{n} = 3$. Panna kirja kõikvõimalikud valimid ja nende saamise tõenäosused järgmisse tabelisse. Poissoni disaini korral anna ise erinevad kaasmistõenäosused π_i objektidele, nii et $\sum_U \pi_i = 3$:

k	LJV TTA	LJV TGA	Bernoulli	Poisson	TP
0000	$p(0000)=?$				
0001					
...					
0003					
...					
3000					
$p(k)$					

2.4 Disaini karakteristikud

Olgu $I \sim p(k)$ valikudisain üldkogumil U . Disainikarakteristikud on arvud, mis kirjeldavad jaotust $p(k)$. Kõige tähtsamad jaotuse $p(k)$ karakteristikud on tema momendid:

$E(I_i)$ – esimest järku moment on objekti i oodatav valikute arv
 $E(I_i I_j)$ – teist järku moment;
 $V(I_i) = E(I_i)^2 - (E I_i)^2$ – valikuindikaatori I_i dispersioon;
 $Cov(I_i, I_j) = \Delta_{ij} = E(I_i I_j) - E(I_i)E(I_j)$ – valikuindikaatorite kovariatsioon
 I_i, I_j ;
 $\Delta_{ij} = 0$ Poissoni disaini korral, muidu $\Delta_{ij} \neq 0$.

Definitsioon. Disaini nimetatakse isekaaluvaks, kui $E(I_i) = const \forall i$.

Definitsioon. Disaini nimetatakse mõõtuvaks, kui $E(I_i I_j) > 0$.

TTA disainide korral:

$I_i \sim B(1, \pi_i)$:

$E(I_i) = P(I_i = 1) = \pi_i$ – esimest järku kaasamistõenäosus;

$E(I_i I_j) = P(I_i = 1, I_j = 1) = \pi_{ij}$ – teist järku kaasamistõenäosus

$\Delta_{ij} = \pi_{ij} - \pi_i \pi_j$

$V(I_i) = \Delta_{ii} = \pi_i - \pi_i^2 = \pi_i(1 - \pi_i)$

Paneme tähele, et $\pi_{ii} = \pi_i$.

Üldjuhul on disaini valimimaht juhuslik suurus, $\mathbf{n} = \sum_U I_i$.

Teoreem Valimimahu \mathbf{n} tähtsamad karakteristikud avalduvad valikudisaini momentide kaudu:

$$E(\mathbf{n}) = \sum_i E(I_i), \quad (1)$$

$$V(\mathbf{n}) = \sum_i \sum_j \Delta_{ij}. \quad (2)$$

Tõestus. Kuna $\mathbf{n} = \sum_{i=1}^N I_i$, siis (1) on tõestatud. Edasi,

$$\begin{aligned} V(\mathbf{n}) &= E \left[\sum_i I_i - \underbrace{E(\sum_i I_i)}_{\sum_i EI_i} \right]^2 = E \left[\sum_i (I_i - EI_i) \right]^2 = \\ &E \left[\sum_i \sum_j (I_i - EI_i)(I_j - EI_j) \right] = \sum_i \sum_j \Delta_{ij}. \quad \diamond \end{aligned}$$

Teoreem Fikseeritud mahuga n disaini $p(k)$ korral kehtivad seosed:

$$\sum_i E(I_i) = n, \quad (3)$$

$$\sum_i \sum_j E(I_i I_j) = n^2, \quad (4)$$

$$\sum_i E(I_i I_j) = nE(I_j), \quad (5)$$

$$\sum_i \sum_j \Delta_{ij} = 0, \quad (6)$$

$$\sum_i \Delta_{ij} = \sum_j \Delta_{ij} = 0. \quad (7)$$

Tõestus. Fikseeritud mahuga disaini korral on $\mathbf{n} = \sum_i I_i \equiv n$ konstant, millest tulenevalt on $E(\mathbf{n}) = n$ ja $V(\mathbf{n}) = 0$, ja seega (3) ning (6) on tõestatud eelmise teoreemi põhjal;

$$(4) \quad \sum_i \sum_j E(I_i I_j) = E \left(\underbrace{\sum_i I_i}_n \underbrace{\sum_j I_j}_n \right) = n^2;$$

$$(5) \quad \sum_i E(I_i I_j) = E \sum_i (I_i I_j) = EI_j \underbrace{\sum_i I_i}_n = nEI_j;$$

$$(7) \quad \sum_i \Delta_{ij} = \sum_i E(I_i I_j) - \sum_i (EI_i)(EI_j) = nEI_j - nEI_j = 0. \quad \diamond$$

3 Hindamise alused

Olgu $U = \{1, 2, \dots, N\}$ üldkogum, ja θ üldkogumi parameeter y_i – uuritav tunnus, mõõdetud objektil $i \in U$.

Oleme huvitatud peamiselt selliste parameetrite θ hindamisest nagu kogusumma $t_y = \sum_U y_i$ või kogusummade suhe $R = \frac{t_y}{t_x}$.

Ka keskmised avalduvad suhetena.

Oleme huvitatud hinnangute $\hat{\theta}$ omadustest, sellistest nagu nihe, dispersioon, dispersiooni hinnang.

Tuletame meelde: Hinnang $\hat{\theta}$ on parameetri θ jaoks nihketa, kui $E\hat{\theta} = \theta$.

Kuna juhuslikkus tekitatakse hinnangusse valikudisaini poolt, siis on keskväärtus defineeritud valikudisaini suhtes:

$$E\hat{\theta} = \sum_k \hat{\theta}(k)p(k),$$











kus summeerimine on üle kõigi võimalike valimite \mathbf{k} . Lähenemist, kus hinnangu keskväärtust, ja sellest tulenevalt ka dispersiooni, defineeritakse valikudisaini abil, nimetatakse disainipõhiseks lähenemiseks. On olemas ka mudelipõhine lähenemine. Mõelgem, kuidas intuiivselt mõista $\hat{\theta}$ disainipõhist keskväärtust ja dispersiooni?

3.1 Valikuuringu andmed

Vaatame esmalt ühte praktilist näidet.

Näide (täiendamiseks loengul). All on toodud kilpkonnade register 10-st kilpkonnast koos nende vanustega (aastates).

Uuritavaks parameetriks on kilbi keskmine paksus (mm). Kirjeldada Poissoni disain $n=4$ valimi võtmiseks.

ÜK	Vanus	π_i	k	w	y	$\tilde{y}_i = w_i \cdot y_i$
1 	100	0,8	1			
2 	90	0,72	0			
3 	70	0,56	1			
4 	55	0,44	0			
5 	50	0,4	1			
6 	40	0,32	0			
7 	40	0,32	0			
8 	35	0,28	1			
9 	15	0,12	0			
10 	5	0,04	0			

Mida teame enne ja mida pärast valiku teostamist ja andmete kogumist?

- Enne valiku teostamist teame:

$U = (1, \dots, N)$ – objektide märgendid (id-kood, nimi, aadress freimis);

$x = (x_1, \dots, x_N)$ – abitunnused, mis on teada iga i jaoks
või mille kogusummad t_x on teada;

$I = (I_1, \dots, I_N) \sim p(k)$ – valikudisain, fikseeritakse planeerimisfaasis,
 $p(k)$ või selle karakteristikud on teada,
 I realisatsioon pole teada.

Ei tea

$y = (y_1, \dots, y_N)$ – uuritava tunnuse väärtusi.

Praktikas on y tegelikult maatriks, milles iga objekti jaoks on veerg paljude tunnuste väärtustega.

$$y_i = \begin{pmatrix} z_i \\ u_i \\ \vdots \end{pmatrix}, \quad i \in U.$$

- Pärast valikut ja andmete kogumist teame:
 $U = (1, \dots, N)$;
 $x = (x_1, \dots, x_N)$;
 $I = (I_1, \dots, I_N)$ realisatsiooni, valimit;
 $y_s = (I_1 y_1, \dots, I_N y_N)$ objektide mõõtmistulemusi valimis, y_i võib olla vektor, nii nagu ka x_i .

Suurused (I, y_s, x) moodustavad valikuuringute andmete stohhastilise esituse. Siin on ilmutatud kujul näha juhuslikkust põhjustav vektor I , $I \sim p(k)$. Selline esitus on aluseks disainipõhisele valikuteooriale. Mudelipõhine valikuteooria eeldab, et $y = (y_1, \dots, y_N)$ ise on juhuslik, juuba üldkogumis. Tema juhuslikku olemust iseloomustatakse mudeliga, näiteks $y_i \sim N(\mu, \sigma^2)$, sõltumatud. Kolmiku (I, y_s, x) funktsiooni nimetame statistikuks.

3.2 Kogusumma nihketa hindamine

Selleks et konstrueerida lihtsaimat hinnangut kogusummale t , vaadeldagem lihtsat statistikut, nimelt andmete

$$y_s = (I_1 y_1, \dots, \underbrace{I_i y_i}_{Y_{si}}, \dots, I_N y_N)$$

lineaarkombinatsiooni

$$\hat{t} = \sum_U c_i y_{si} = \sum_U c_i I_i y_i,$$

kus c_i on mittejuhuslik konstant y_{si} on objekti i vaatlustulemus.

Tahame, et

$$E\left[\sum_U c_i I_i y_i\right] = \sum_U c_i y_i E(I_i) = \sum_U y_i \quad - \text{nihketus.}$$

Et seos kehtiks, peab olema

$$c_i = \frac{1}{E(I_i)}.$$

Järelikult saame kogusumma nihketa hinnaguks

$$\hat{t} = \sum_U \frac{I_i y_i}{E(I_i)}. \quad (8)$$

Valikudisainide nõue $E(I_i) > 0$ on nüüd selge. Antud hinnangul on kaks tähtsat esitust:

$$\hat{t} = \sum_U I_i \check{y}_i, \quad (9)$$

kus

$$\check{y}_i = \frac{y_i}{E(I_i)} - \text{laiendatud } y_i,$$

ja teiseks,

$$\hat{t} = \sum_U \omega_i y_i, \quad (10)$$

kus

$$\omega_i = \frac{I_i}{E(I_i)} - y_i \text{ valikukaal.}$$

Pangem tähele, et $\omega_i = 0$ mittevalitud objektide jaoks ja ta kaalub üles valitud objekte. Valimiväärtus y_i esindab ω_i objekti üldkogumis (tavaliselt $\omega_i \gg 1$).

Näeme valitud objektide erinevat panust hinnangusse. Need, mille valik on oodatavalt suurem ($E(I_i)$ suurem), surutakse maha väiksema kaaluga ω_i .

Pange tähele ja pidage meeles, et kuigi summeerimine toimub üle U , esitavad valemid (8)-(10) ikkagi valimisummasid, ja seega on arvutatavad andmetelt. Realiseerunud valimi korral on need summad tegelikult üle valimi s :

$$\hat{t} = \sum_s I_i \check{y}_i \quad \text{or} \quad \hat{t} = \sum_s \omega_i y_i,$$

kus I_i ja ω_i on realiseerunud väärtused ja s loendab üldkogumist U valitud erinevaid objekte:

$$s = \{i : i \in U, \text{ mille korral } I_i > 0\}.$$

Teoreetilises käsitluses eelistame selles konseptsis summasid üle U . Hinnangu (8) dispersiooniavaldise saame kasutades tuntud seost tõenäosusteooriast:

$$V\left[\sum_{i=1}^N c_i X_i\right] = \sum_{i=1}^N \sum_{j=1}^N c_i c_j \text{Cov}(X_i, X_j), \quad X_i - \text{juhuslik suurus.} \quad (11)$$

Seoses on lihtne veenduda, kasutades definitsioone $V(X) = E(X - EX)^2$, $\text{Cov}(X, Y) = E(X - EX)(Y - EY)$.

Rakendades seost (11) meie hinnangule $\hat{t} = \sum_U I_i \check{y}_i$, saame

$$V(\hat{t}) = \sum_{i=1}^N \sum_{j=1}^N \Delta_{ij} \check{y}_i \check{y}_j, \quad (12)$$

kus $\Delta_{ij} = \text{Cov}(I_i, I_j)$. Dispersioon (12) on hinnangu \hat{t} disainipõhine dispersioon. Ta on hinnangu \hat{t} varieeruvuse mõõt antud valikudisaini $p(k)$ korral. $V(\hat{t})$ valemis (12) on teoreetiline avaldis, ta ei ole arvutatav.

Vajame suuruse $V(\hat{t})$ hinnangut.

Hinnangus ei saa kasutada üldkogumiväärtusi \check{y}_i . Saab kasutada valimiväärtusi $I_i \check{y}_i$.

Kirjutagem (12) valimiväärtuste kaudu ja lisagem tundmatud konstandid c_{ij} :

$$\hat{V}(\hat{t}) = \sum_{i=1}^N \sum_{j=1}^N c_{ij} \Delta_{ij} I_i \check{y}_i I_j \check{y}_j.$$

Konstandid c_{ij} määrame nihketuse nõudest:

$$E\left[\hat{V}(\hat{t})\right] = V(\hat{t}), \quad (13)$$

$$E\left[\hat{V}(\hat{t})\right] = \sum \sum_U c_{ij} \Delta_{ij} \check{y}_i \check{y}_j E(I_i I_j).$$

On selge, et tingimus (13) kehtib, kui valime $c_{ij} = \frac{1}{E(I_i I_j)}$.

Seetõttu on $V(\hat{t})$ nihketa hinnanguks

$$\hat{V}(\hat{t}) = \sum \sum_U \check{\Delta}_{ij} \check{y}_i \check{y}_j I_i I_j,$$

kus $\check{\Delta}_{ij} = \frac{\Delta_{ij}}{E(I_i I_j)}$ on valikumuutujate laiendatud kovariatsioon. Dispersiooni hinnangu saame alternatiivselt esitada valikukaalude abil

$$\hat{V}(\hat{t}) = \sum \sum_U \check{\Delta}_{ij} \omega_i y_i \omega_j y_j.$$

Märgime veelkord: kui summa üle U sisaldab valikumuutujaid I_i või valikukaalu ω_i , toimub summeerimine tegelikult üle valimi s , ja seega valem esitab hinnangut.

Saadud tähtsad tulemused on koondatud järgmisse teoreemi.

Teoreem(Üldine hindamisteoreem) Üldkogumi kogusumma $t = \sum_U y_i$ nihketa hinnang on

$$\hat{t} = \sum_U I_i \check{y}_i \quad (\text{või } \hat{t} = \sum_U \omega_i y_i), \quad (14)$$

kus

$$\check{y}_i = \frac{y_i}{E(I_i)} \quad \text{ja} \quad \omega_i = \frac{I_i}{E(I_i)}. \quad (15)$$

Selle disainipõhine dispersioon on

$$V(\hat{t}) = \sum \sum_U \Delta_{ij} \check{y}_i \check{y}_j, \quad (16)$$

kus $\Delta_{ij} = Cov(I_i, I_j)$. Dispersiooni nihketa hinnanguks $E(I_i I_j) > 0$ korral on

$$\hat{V}(\hat{t}) = \sum \sum_U \check{\Delta}_{ij} \check{y}_i \check{y}_j I_i I_j \quad (\text{või } \hat{V}(\hat{t}) = \sum \sum_U \check{\Delta}_{ij} \omega_i y_i \omega_j y_j), \quad (17)$$

kus

$$\check{\Delta}_{ij} = \frac{\Delta_{ij}}{E(I_i I_j)}.$$

Märkus Üldine hindamisteoreem kehtib iga valikudisaini korral, nii TTA kui TGA disainide korral. Vaja on vaid teada disainikarakteristikuid

$$E(I_j), \quad E(I_i I_j), \quad \Delta_{ij} \quad \text{for } i = j \text{ ja } i \neq j.$$

Maatriksite abil saab kahekordsed summad sageli elegantsemalt esitada. Nii saame dispersiooni jaoks avaldise:

$$V(\hat{t}) = \check{y}' \Delta \check{y}, \quad (18)$$

kus $\Delta = (\Delta_{ij}) : N \times N$ and $\check{y} = (\check{y}_i) : N \times 1$. Sarnaselt saame dispersiooni hinnangu jaoks

$$\hat{V}(\hat{t}) = \check{y}'_s \check{\Delta} \check{y}_s, \quad (19)$$

kus $\check{\Delta} = (\check{\Delta}_{ij}) : N \times N$ ja $\check{y}_s = (\check{y}_i I_i) : N \times 1$. Kuna mittevalitud elemendid vektoris \check{y}_s on nullid, siis saavutavad ruutvormi (19) komponendid väiksema dimensiooni $\Delta : n \times n$ ja $\check{y}_s : n \times 1$, kus n on valimimaht.

Seda maatriksesitust kasutame hiljem IML programmis.

3.2.1 Dispersiooni hindamine fikseeritud mahuga disainide korral

Tuletame meelde, et hinnangu dispersioon on tema hajuvuse mõõt. Kui valikudisain on fikseeritud, siis hinnangu dispersioon on teatav arvuline konstant (tavaliselt küll tundmatu). Samas selle ühe konstandi jaoks saab konstrueerida mitmeid hinnanguid. Lisaks eespool toodud üldisele hinnangule vaatame siin teist hinnangut, mis kehtib üksnes fikseeritud mahuga disainide korral.

Olgu disain $p(k)$ fikseeritud valimimahuga $-\sum_U I_i \equiv n$.

Teoreem Fikseeritud mahuga disaini $p(k)$ korral saab hinnangu $\hat{t} = \sum_U I_i \check{y}_i$ dispersiooni esitada alternatiivsel kujul

$$V(\hat{t}) = -\frac{1}{2} \sum \sum_U \Delta_{ij} (\check{y}_i - \check{y}_j)^2, \quad (20)$$

ja eeldusel, et $E(I_i I_j) > 0 \forall i \neq j \in U$, on dispersiooni $V(\hat{t})$ nihketa hinnanguks

$$\hat{V}(\hat{t}) = -\frac{1}{2} \sum \sum_U I_i I_j \check{\Delta}_{ij} (\check{y}_i - \check{y}_j)^2. \quad (21)$$

Tõestus. Näitame, et (20) on ekvivalentne seosele (16). Avame sulud seoses (20):

$$V(\hat{t}) = -\frac{1}{2} \sum_i \sum_j \Delta_{ij} (\check{y}_i^2 - 2\check{y}_i \check{y}_j + \check{y}_j^2).$$

$$\text{Nüüd} \quad \sum_i \sum_j \Delta_{ij} \check{y}_i^2 = \sum_i \check{y}_i^2 \underbrace{\sum_j \Delta_{ij}}_0 = 0.$$

$$\text{Samuti} \quad \sum_i \sum_j \Delta_{ij} \check{y}_i^2 = 0, \text{ ja saame}$$

$$V(\hat{t}) = \sum_{i=1}^N \sum_{j=1}^N \Delta_{ij} \check{y}_i \check{y}_j,$$

mis on (16).

On lihtne näha, et (21) on nihketa suuruse (20) jaoks. \diamond

Märkus Avaldis (16) annab hinnangu \hat{t} dispersiooni kõigi valikudisainide jaoks. Avaldis (20) annab selle üksnes fikseeritud mahuga disainide jaoks, millisel juhul ta on võrdne seoses (20) antuga. Aga NB! dispersiooni hinnang (21) ei ole üldjuhul võrdne hinnanguga (17), isegi mitte fikseeritud mahuga disainide korral.

Kui valikudisain on fikseeritud mahuga, siis eelistatakse dispersioonihinnangut (21), kuna see on stabiilsem (tal on väiksem varieeruvus üle erinevate valimite) ja üldjuhul ei tule ta negatiivne ($\Delta_{ij} < 0$ enamuse praktikas kasutatavate fikseeritud mahuga disainide korral).

Dispersioonihinnangut (21) nimetatakse Sen–Yates–Grundy hinnanguks.

Jätkame 3. loengu materjaliga... SYG teoreem...

Näide (täiendamiseks loengul). Jätkame kilpkonnade andmetega. Olgu tegemist LJV TTA (fikseeritud mahuga disain). $N = 10$, $n = 4$. Sel juhul $EI_i = np_i = n/N = 0,4$.

Leiame keskmise kilbi paksuse hinnangu, hinnangu dispersioon ja dispersiooni hinnangu nii ÜHT kui SYG valemite järgi.

ÜK	Vanus	EI_i	\mathbf{k}	\mathbf{w}	\mathbf{y}	$\tilde{y}_i = w_i \cdot y_i$
1	100	0,4	1	2,5	2,3	5,75
2	90	0,4	0	0	0	0
3	70	0,4	2	5	2,1	5 * 2,1 = 10,5
4	55	0,4	0	0	0	0
5	50	0,4	0	0	0	0
6	40	0,4	0	0	0	0
7	40	0,4	0	0	0	0
8	35	0,4	1	2,5	0,7	1,75
9	15	0,4	0	0	0	0
10	5	0,4	0	0	0	0

3.2.2 Nihketa hinnang kogusummale TGA disainide korral

Vaatleme tähtsaimat ebavõrdsete tõenäosustega TGA disaini – multinomiaaldisaini:

$$I \sim M(n; p_1, p_2, \dots, p_N), \sum_U p_i = 1, I_i \sim B(n, p_i).$$

Sel juhul

$$\begin{aligned} E(I_i) &= np_i; \\ \Delta_{ii} &= V(I_i) = np_i(1 - p_i); \\ \Delta_{ij} &= Cov(I_i, I_j) = -np_i p_j; \\ E(I_i I_j) &= n(n-1)p_i p_j; \text{ Näidata!} \\ E(I_i^2) &= np_i(1 - p_i + np_i); \text{ Näidata!} \\ w_i &= \frac{I_i}{np_i}. \end{aligned}$$

Hinnangufunktsioon saab järgmist kuju:

$$\hat{t} = \sum_U \frac{I_i y_i}{np_i}.$$

Seda tüüpi hinnangut nimetatakse Hansen- Hurwitz hinnanguks ja ka p-hinnanguks.

Teoreem: hindamisteoreem multinomiaaldisaini korral Multinomiaaldisaini korral on nihketa hinnang ÜK kogusummale $t = \sum_U y_i$ järgmine:

$$\hat{t} = \sum_U \frac{I_i y_i}{np_i}.$$

Hinnangu \hat{t} dispersioon on:

$$V(\hat{t}) = \frac{1}{n} \left[\sum_U \frac{y_i^2}{p_i} - t^2 \right].$$

Dispersiooni kaks nihketa hinnangut on:

$$(1) \hat{V}(\hat{t}) = \frac{1}{n-1} \left[\sum_U \frac{n}{1-p_i+np_i} \left(\frac{I_i y_i}{np_i} \right)^2 - \hat{t}^2 \right],$$

$$(2) \hat{V}(\hat{t}) = \frac{1}{n(n-1)} \left[\sum_U I_i \frac{y_i^2}{p_i^2} - n\hat{t}^2 \right].$$

Märkus. Dispersioonihinnang (1) jäeldub Üldisest hindamisteoreemist. Dispersioonihinnang (2) jäeldub Teoreemist fikseeritud mahuga disainide kohta. Viimast hinnangut nimetatakse Sen-Yates-Grundy (SYG) dispersioonihinnanguks. Ta on lihtsama kujuga ja ta on ka stabiilsem kui dispersioonihinnang (1).

Tõestus.

Märkus Juhul, kui multinomiaaldisaini korral

$$y_i = cp_i, \quad i = 1, \dots, N,$$

siis ka $V(\hat{t}) = 0$ (näita!). St et kui uuritava tunnuse väärtused on võrdelised tõenäosustega p_i , siis hinnang \hat{t} annab täpse ÜK summa.

Praktikas pole võimalik valimi võtmisel kasutada selliseid tõenäosusi p_i , seda enam, et uuritavaid tunnuseid on palju. Kui osatakse määrata sellised p_i , mis on ligikaudu võrdelised väärtustega y_i , siis saavutatakse väiksem dispersioon vastavale kogusumma hinnangule.

Üritatakse leida selline tausttunnust x , mis on teada kõigi üldkogumi objektide kohta ja mis on positiivselt tugevasti korreleeritud uuritava tunnusega y . Selle abil määratakse

$$p_i = \frac{x_i}{\sum_U x_i}; \quad i = 1, \dots, N.$$

Valikuuringutes on selliseks tausttunnuseks sageli objekti suuruse tunnus, mistõttu multinomiaaldisaini nimetatakse ka suurusega võrdeliste tõenäosusega disainiks (*pps* ehk *probability proportional-to-size sampling*).

3.2.3 Üldkogumi keskmise nihketa hindamine

ÜK keskmine, täpsemalt keskmine objekti kohta, on defineeritud järgmiselt:

$$\bar{Y} = \frac{1}{N} \sum_U y_i = \frac{t_y}{N}.$$

Kui N on teada, saab sellele anda nihketa hinnangu.

(1) N on teada => piisab kogusumma hindamisest:

$$\hat{Y} = \frac{\hat{t}_y}{N}. \quad (*)$$

Dispersioon ja dispersioonihinnang järelduvad teadaolevates tulemustest \hat{t}_y kohta:

$$\begin{aligned} V(\hat{Y}) &= \frac{1}{N^2} V(\hat{t}_y), \\ \hat{V}(\hat{Y}) &= \frac{1}{N^2} \hat{V}(\hat{t}_y). \end{aligned}$$

(2) N pole teada \Rightarrow saab kasutada alternatiivset hinnangut:

$$\hat{Y}_{alt} = \frac{\hat{t}_y}{\hat{N}}, \text{ kus } \hat{N} = \sum_U w_i, w_i = \frac{I_i}{E(I_i)}. \quad (**)$$

Paneme tähele, et nimetajas on disainikaalude summa, mis on nihketa hinnanguks üldkogumi mahule N . Hinnangu \hat{N} omadused tulenevad Üldisest hinadamisteoreemist erijuhul, kui $y_i \equiv 1$. Siis $N = \sum_U 1$ ja $\hat{N} = \sum_U w_i 1$.

Keskmise alternatiivse hinnangu dispersiooni kohe leida ei saa, sest tegemist on kahe juhusliku suuruse suhtega. Dispersioonivalemiteni jõuame hiljem, kui vaatame suhte hindamist üldjuhul.

Märkus 1. Isegi kui N on teada, eelistatakse hinnangut $(**)$ hinnangule $(*)$, kuna üldjuhul on $(**)$ väiksema varieeruvusega. Hinnang $(**)$ annab väikese nihkega tulemuse, kuid see nihe on vähe oluline võrreldes dispersiooniga.

Märkus 2. Ka kogusumma \hat{t}_y hindamisel eelistatakse sageli järgmist hinnangut, mida nimetatakse kogusumma suhtehinnanguks ja mis on regressioonhinnangu erijuht:

$$\hat{t}_{y,alt} = \hat{Y}_{alt} \cdot N = \frac{N}{\hat{N}} \hat{t}_y.$$

Näeme, et see kogusumma hinnang vajab lisainformatsiooni ehk siin N teadmist.

Näide. Olgu tegemist multinomiaaldisainiga, kilpkonnade andmestik mahuga $N = 100$. Soovime valimit mahuga $n = 40$, kus $p_i = \frac{V_{anus_i}}{\sum_U V_{anus_i}}$.

Uuritavaks tunnuseks on keskmine kilbi paksus, $\bar{Y} = \frac{t}{N}$. Sellele pakkume kahte hinnangut, $\hat{Y} = \frac{\hat{t}}{N}$ ja $\hat{Y}_{alt} = \frac{\hat{t}}{\hat{N}}$ ning uurime hinnangute varieeruvust simuleerimise abil.

Simulatsioon (kõik y tunnuse väärtused on ÜK-s teada):

1. samm. Võtame pps-valimi mahuga 40 ja leiame \hat{Y} ning \hat{Y}_{alt} .
2. samm. Kordame 1.sammu 1000 korda. Hinnanguid salvestame vastavatesse vektoritesse.

3. samm. Joonestame mõlema hinnangu histogrammi ning võrdleme jaotused visuaalselt.

Korduvate valimite võtmiseks sobib selline SAS programm.

```
proc surveyselect
  data= kilpkonnad
  sampsize=4|
  method=pps
  reps=1000
  ;
  size vanus;
run;
```

Tulemusena on andmestik, milles on 1000 valimit mahuga 40. Neid andmeid kasutades leiame hinnangud näiteks SQL abil.

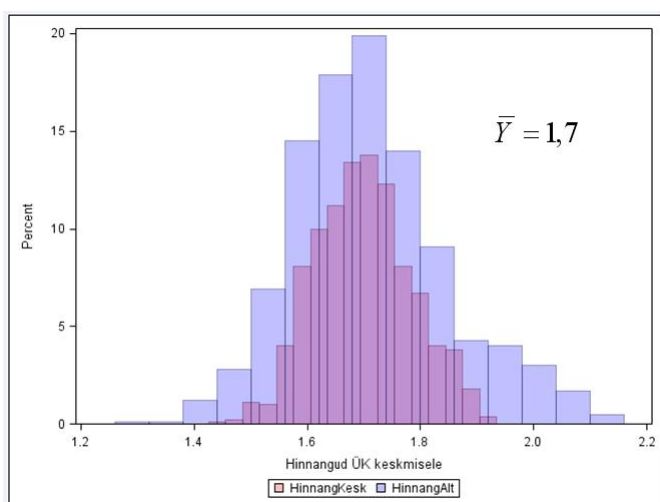
```
proc sql;
  create table Hinnangud as
  select sum(y*SamplingWeight)/10
         as HinnangKesk,
         sum(y*SamplingWeight)/sum(SamplingWeight)
         as HinnangAlt
  from Valimid
  group by Replicate;
quit;
```

Mõlema hinnangu jaotust saab ühele diagrammile lisada alljärgneva protseduuri abil:

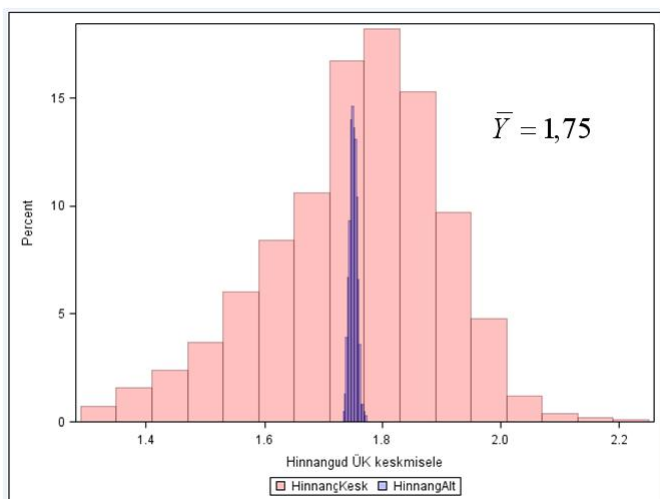
```
proc sgplot data=hinnangud;
  histogram HinnangKesk /
    transparency=0.75 fillattrs=(color=red);
  histogram HinnangAlt /
    transparency=0.75 fillattrs=(color=blue);
  keylegend / location=outside position=bottom;
  axis label="Hinnangud ÜK keskmisele";
run;
```

Osutub, et tulemus sõltub disainikaaludest, ehk valikutõenäosustest p_i . Ju-

hul, kui uuritav tunnus y on tugevasti ja positiivselt korrileeritud registritunnusega $Vanus$ (mille järgi valikutõenäosused olid valitud), siis hinnangute varieeruvus väga ei erine. Isegi tundub tavahinnangu \hat{Y} varieeruvus väiksem.



Kui aga uuritava tunnuse ja registri tunnuse vahel seost pole või see on isegi negatiivne, siis alternatiivne hinnang osutub mitu korda paremaks.



Selleks, et oleks võimalik uurida hinnangute dispersioone teoreetiliselt, läheb vaja veel ühte mõistet - hinnangutevahelist kovariatsiooni.

3.2.4 Kahe nihketa hinnangu kovariatsioon

Vaatame kahte kogusummat t_y ja t_x . Olgu vastavad nihketa hinnangud

$$\hat{t}_y = \sum_U I_i \check{y}_i \text{ ja } \hat{t}_x = \sum_U I_i \check{x}_i,$$

kus $\check{y}_i = y_i/EI_i$ ja $\check{x}_i = x_i/EI_i$. Siis kovariatsiooni definitsiooni järgi:

$$Cov(\hat{t}_y, \hat{t}_x) = E [(\hat{t}_y - t_y)(\hat{t}_x - t_x)].$$

Paneme tähele, et

$$\hat{t}_y - t_y = \sum_U I_i \check{y}_i - \sum_U y_i = \sum_U I_i \check{y}_i - \sum_U E(I_i) \check{y}_i = \sum_U [I_i - E(I_i)] \check{y}_i,$$

samuti

$$\hat{t}_x - t_x = \sum_U [I_i - E(I_i)] \check{x}_i.$$

Kasutades seost

$$\left(\sum_{i=1}^N a_i \right) \left(\sum_{i=1}^N b_i \right) = \sum_{i=1}^N \sum_{j=1}^N a_i b_j$$

saame,

$$Cov(\hat{t}_y, \hat{t}_x) = E \left[\sum_i \sum_j (I_i - EI_i) \check{y}_i (I_j - EI_j) \check{x}_j \right] = \sum_i \sum_j \underbrace{E(I_i - EI_i)(I_j - EI_j)}_{Cov(I_i, I_j) = \Delta_{ij}} \check{y}_i \check{x}_j.$$

Järelikult,

$$Cov(\hat{t}_y, \hat{t}_x) = \sum_{i=1}^N \sum_{j=1}^N \Delta_{ij} \check{y}_i \check{x}_j$$

nihketa hinnanguga:

$$\hat{C}ov(\hat{t}_y, \hat{t}_x) = \sum_{i=1}^N \sum_{j=1}^N \check{\Delta}_{ij} \check{y}_i \check{x}_j I_i I_j,$$

kus endiselt $\check{\Delta}_{ij} = \Delta_{ij}/E(I_i I_j)$.

3.3 Suhte hinnang

Olgu uuritavaks parameetrike kahe tunnuse kogusummade jagatis

$$R = \frac{t_y}{t_x}.$$

Näiteks soovime uurida pere kulutuste osakaalu meelelahutusele. Parameetri R hinnangu leidmiseks peame hindama t_y ja t_x :

$$\hat{R} = \frac{\hat{t}_y}{\hat{t}_x}.$$

See hinnang on kahe juhusliku suuruse mittelineaarne funktsioon ja tema statistiliste omaduste täpne tuletamine ei ole lihtne.

Ligikaudseks tuletamiseks kasutatakse lineariseerimistehnikat Tayloriga reaktiivide arenduse abil.

3.3.1 Tayloriga rida kahe argumendi korral

Olgu X_1, X_2 juhuslikud suurused ja $g(X_1, X_2)$ nende funktsioon. Funktsiooni $g(\cdot, \cdot)$ Tayloriga rea lineaarosa punkti (a_1, a_2) ümbruses on järgmine:

$$g(X_1, X_2) \approx g(a_1, a_2) + \left. \frac{\partial g}{\partial X_1} \right|_{(a_1, a_2)} (X_1 - a_1) + \left. \frac{\partial g}{\partial X_2} \right|_{(a_1, a_2)} (X_2 - a_2),$$

kus $\left. \frac{\partial g}{\partial X_1} \right|_{(a_1, a_2)}$ on g osatuletis punktis (a_1, a_2) .

Saadud osatuletised pole enam juhuslikud suurused ja seega saadud avaldis on lineaarne funktsioon X_1 ja X_2 suhtes.

3.3.2 Suhte hinnangu Tayloriga rittaarendus

Arendame $\hat{R} = \frac{\hat{t}_y}{\hat{t}_x} = g(\hat{t}_y, \hat{t}_x)$ Tayloriga ritta punkti (t_y, t_x) ümbruses.

$$\left. \frac{\partial \hat{R}}{\partial \hat{t}_y} \right|_{(t_y, t_x)} = \frac{1}{t_x},$$

$$\left. \frac{\partial \hat{R}}{\partial \hat{t}_x} \right|_{(t_y, t_x)} = -\frac{t_y}{t_x^2}.$$

Seega,

$$\hat{R} \approx \frac{t_y}{t_x} + \frac{1}{t_x}(\hat{t}_y - t_y) - \frac{t_y}{t_x^2}(\hat{t}_x - t_x) = R + \frac{1}{t_x}(\hat{t}_y - R\hat{t}_x).$$

Veendume, et saadud \hat{R} on ligikaudu nihketa:

$$E(\hat{R}) \approx R + \frac{1}{t_x} \underbrace{[E(\hat{t}_y)]}_{t_y} - R \underbrace{[E(\hat{t}_x)]}_{t_x} = R.$$

Leiame dispersiooni:

$$V(\hat{R}) = \frac{1}{t_x^2} V(\hat{t}_y - R\hat{t}_x) = \frac{1}{t_x^2} [V(\hat{t}_y) + R^2 V(\hat{t}_x) - 2RCov(\hat{t}_y, \hat{t}_x)],$$

kus kasutame lineaarkombinatsiooni dispersioonivalemit:

$$V(aX + bY) = a^2 V(X) + b^2 V(Y) + 2ab Cov(X, Y).$$

Selleks, et saada dispersioonile $V(\hat{R})$ hinnangut, võime kasutada nihketa hinnanguid \hat{t}_x , \hat{t}_y ja $\hat{V}(\hat{t}_x)$, $\hat{V}(\hat{t}_y)$ üldisest hindamisteoreemist. Teame ka kovariatsiooni nihketa hinnangut:

$$\hat{V}(\hat{R}) = \frac{1}{\hat{t}_x^2} \left[\hat{V}(\hat{t}_y) + \hat{R}^2 \hat{V}(\hat{t}_x) - 2\hat{R}\hat{C}ov(\hat{t}_y, \hat{t}_x) \right].$$

Suhte hinnangu dispersiooni leidmiseks on olemas ka alternatiivne valem. Alternatiivne valem võimaldab suhte hinnangut käsitleda varasema teooria valguses (ÜHT jm tulemused). Vajalikuks osutub sobiva uue tunnuse defineerimine.

Alternatiivseks esituseks kirjutame Tayloriga lineaarliikme veel kord välja:

$$\hat{R} \approx R + \frac{1}{t_x}(\hat{t}_y - R\hat{t}_x) = R + \frac{1}{t_x} \left[\sum_U I_i \check{y}_i - R \sum_U I_i \check{x}_i \right] = R + \frac{1}{t_x} \sum_U I_i (\check{y}_i - R\check{x}_i). \quad (22)$$

Viimane saadud summa on nihketa hinnang kogusummale

$$\sum_U (y_i - Rx_i)$$

ja me saame kasutada ÜHT dispersioonivalemi saamiseks. Võttes kasutusele uue tunnuse

$$u_i = y_i - Rx_i, \quad (23)$$

saame suhte hinnangu esitada kujul

$$\hat{R} \approx R + \sum_U I_i \check{u}_i,$$

millest saame

$$V(\hat{R}) \approx \frac{1}{t_x^2} \sum_{i=1}^N \sum_{j=1}^N \Delta_{ij} \check{u}_i \check{u}_j. \quad (24)$$

Dispersiooni hinnang on vastavalt

$$\hat{V}(\hat{R}) = \frac{1}{\hat{t}_x^2} \sum_{i=1}^N \sum_{j=1}^N \check{\Delta}_{ij} \check{u}_i \check{u}_j I_i I_j, \quad (25)$$

kus

$$\hat{u}_i = y_i - \hat{R}x_i, \quad (26)$$

$$\check{u}_i = \check{y}_i - \hat{R}\check{x}_i. \quad (27)$$

Kaalude abil on dispersioonihinnangu valem:

$$\hat{V}(\hat{R}) = \frac{1}{\hat{t}_x^2} \sum_{i=1}^N \sum_{j=1}^N \check{\Delta}_{ij} w_i \hat{u}_i w_j \hat{u}_j.$$

Paneme tähele, et $\hat{t}_y - R\hat{t}_x$ valemis (22) hindab nulli, kuna $0 = t_y - \frac{t_y}{t_x}t_x = t_y - Rt_x$. Järelikult hinnang $\hat{t}_y - R\hat{t}_x$ varieerub nulli ümbruses. Kui me jagame selle hinnangu t_x -ga, siis hinnangu varieeruvus muutub veelgi väiksemaks. Seega on suhte R hinnangu dispersioon väike, mistõttu on \hat{R} hea hinnang.

3.4 Osakogumi hindamine.

Osakogumiks nimetatakse üldkogumi U alamhulka U_d , $U_d \subset U$. Osakogumi mahtu tähistatakse N_d . Osakogumi objektid on sama tüüpi nagu üldkogumi

omad. Näiteks kui üldkogumi objektideks on pered, siis ka osakogum on teatava tunnuse alusel määratud pered (mitte isikud, lapsed vms.)

Mõned Osakogumi näited:

1. Välisosalusega ettevõtted kõigi ettevõtete hulgas Eestis.
2. Töötud riigi tööealiste elanike hulgas.
3. Suitsetajad kopsuvähi haigete hulgas.

Osakogumid määratakse mingi tunnuse (osakogumi identifikaator) järgi. Näiteks tunnus suitsetamine, tunnus tõine staatus jt.

Tavaliselt ei võeta valimit eraldi igas huvipakkavas osakogumis, vaid see võetakse üldkogumis tervikuna. Nii võib juhtuda, et kui osakogumi maht N_d on väike, siis osakogumist U_d satub valimisse vähe objekte ja arvutatavad osakogumi hinnangud on väga väikese täpsusega.

Osakogumit nimetatakse väikeseks osakogumiks, kui valimimaht temas on väike (isegi 0). Selliste osakogumite jaoks on omad hindamismeetodid (Small area estimation methods). Need kasutavad mudeleid, et kompenseerida valimi väiksust. Neid meetodeid me valikuteooria baaskursuses ei vaata. Meie vaatame osakogumite disainipõhist hindamist, mis eeldab, et valimimahud osakogumites ei ole väga väikesed.

Valimimaht ja osakogumi valimimaht avalduvad valikuindikaatorite abil järgmiste summadena,

$$\mathbf{n} = \sum_U I_i, \quad \mathbf{n}_d = \sum_{U_d} I_i.$$

Näeme, et isegi kui \mathbf{n} on fikseeritud, jääb \mathbf{n}_d juhuslikuks.

Huvi pakuvad järgmised parameetrid:

$$\begin{aligned} N_d & & - & \text{osakogumi maht,} \\ P_d = \frac{N_d}{N} & & - & \text{osakogumi osakaal,} \\ t_d = \sum_{U_d} y_i & & - & \text{osakogumi kogusumma,} \\ \bar{Y}_d = \frac{1}{N_d} t_d & & - & \text{osakogumi keskmine,} \\ R_d = \frac{\sum_{U_d} y_i}{\sum_{U_d} x_i} & & - & \text{suhe osakogumis.} \end{aligned}$$

Võtame kasutusele indikaatortunnuse z , mis näitab kuuluvust osakogumisse:

$$z_i = \begin{cases} 1, & i \in U_d, \\ 0, & \text{vastasel juhul.} \end{cases}$$

Nüüd saame osakogumi mahu kirja panna ÜK summana,

$$N_d = \sum_U z_i,$$

mistõttu saame rakendada üldist hindamisteoreemi ÜK summade hindamiseks:

$$\hat{N}_d = \sum_U I_i \check{z}_i, \text{ kus } \check{z}_i = z_i/E(I_i),$$

ehk

$$\hat{N}_d = \sum_U w_i z_i, \text{ kus } w_i = I_i/E(I_i).$$

Ka \hat{N}_d dispersioon ja dispersioonihinnang tulevad üldisest hindamisteoreemist. Tunnus y_i tuleb vaid asendada indikaatortunnusega z_i .

Hinnang \hat{N}_d , kuigi kirja pandud üldkogumi summana, on tegelikult valimisumma. Tähistame valimi (hulkvalimi) seda osa, mis kuulub osakogumisse s_d , st

$$s_d = s \cap U_d.$$

Näeme, et nendes tähistustes

$$\hat{N}_d = \sum_U w_i z_i = \sum_s w_i z_i = \sum_{s_d} w_i,$$

ja seega hinnatud osakogumi maht on kaalude summa üle osavalimi s_d .

Osakogumi y -tunnuse summa $t_d = \sum_{U_d} y_i$ hindamiseks loome uue tunnuse y' :

$$y'_i = z_i y_i = \begin{cases} y_i, & i \in U_d, \\ 0, & \text{vastasel juhul.} \end{cases}$$

Nüüd saame osakogumi summa kirja panna uue tunnuse ÜK summana:

$$t_d = \sum_U y'_i.$$

Selle summa nihketa hinnang

$$\hat{t}_d = \sum_U w_i y'_i$$

ja tema dispersiooni hinnang tulevad jälle üldisest hindamisteoreemist. Valemities tuleb kasutada vaid uut tunnust y'_i .

Kui teame osakogumi mahtu N_d , siis osakogumi keskmise hinnang on lihtsalt

$$\hat{Y}_d = \frac{1}{N_d} \hat{t}_d,$$

mille dispersioon tuleb sellest, et teame \hat{t}_d dispersiooni (ÜHT).

Nägime, et osakogumi mahu ja summa hindamisel saame kasutada üldist hindamisteoreemi.

Kui aga osakogumimaht pole teada, siis saab keskmise hinnang järgmise kuju

$$\hat{Y}_d = \frac{\hat{t}_d}{\hat{N}_d}, \quad (28)$$

mis on kahe hinnangu suhe (suhte hinnang) ja seda tüüpi hinnangu dispersiooni leidmist vaatasime eelmises punktis.

Osakogumite korral huvitatakse ka kahe summa suhtest üldisemal kujul,

$$R_d = \frac{\sum_{U_d} y_i}{\sum_{U_d} x_i},$$

mille saab jällegi esitada üldkogumi U tasemel

$$R_d = \frac{\sum_U y_i z_i}{\sum_U x_i z_i} = \frac{\sum_U y'_i}{\sum_U x'_i}.$$

Nüüd saame nii R_d hinnangu kui ka dispersioonihinnangu taandada juba olemasolevatele valemitele. Näiteks,

$$\hat{R}_d = \frac{\sum_U w_i y'_i}{\sum_U w_i x'_i}, \text{ kus } w_i = I_i/E(I_i)$$

Märkus. Isegi kui teame osakogumi mahtu N_d , on soovitatav kasutada osakogumi keskmise hindamisel hinnangut (28), kuna see on väiksema varieeruvusega. Sellest järelduvalt on kogusumma hindamiseks soovitatav kasutada

$$\hat{t}_d = N_d \hat{Y}_d = N_d \frac{\hat{t}_d}{\hat{N}_d}.$$

3.5 Hinnangu täpsuse iseloomustamine

Olgu θ meid huvitav parameeter ja $\hat{\theta}$ on selle nihketa hinnang.

$\hat{\theta} - \theta$ on viga, mille võivad põhjustada järgmised komponendid:

- **valikuviga** (*sampling error*): valimi juhuslikusest põhjustatud viga; seda on võimalik hinnata, kui on teada valikudisain $p(k)$ või tema karakteristikud, ja hinnangu $\hat{\theta}$ avaldis
- **muu viga**: kaost/mittevastamisest põhjustatud viga, inervjueeriast põhjustatud viga, valesti sõnastatud ankeedist.... Seda laadi viga on raske hinnata, kuid on võimalik määrata vea suundumust (üle/alahinnang).

Uuringu üheks kvaliteedinäitajaks on vastamismäär, $\frac{\text{vastanute arv}}{n}$. Eesti Statistikaameti leibkonnauuringud (*Household Budget Surveys*) toimuvad regulaarselt kord kuus alates 1995. aastast. Vastamismäär on ligikaudu 50%. Eri-nevates riikides erinevate uuringute vastamismäär võib olla 20 – 80%. Kaost põhjustatud nihke vähendamiseks kasutatakse tänapäeval mitmesuguseid kalibreerimismeetodeid. Need vajavad lisainformatsiooni oma konstruktsioonis. Kalibreerimishinnanguid antud kursuses ei vaata.

Valikuviga iseloomustavad järgmised suurused:

- (1) $\sqrt{\hat{V}(\hat{\theta})} - \hat{\theta}$ standardhälbe hinnang, standardviga (standard error);
- (2) $\lambda_{\alpha/2} \sqrt{\hat{V}(\hat{\theta})}$ - pool usaldusvahemiku pikkust;
- (3) $CV(\hat{\theta}) = \frac{\sqrt{\hat{V}(\hat{\theta})}}{\hat{\theta}}$ - suhteline viga (relative error, coefficient of variation)

- (4) $\frac{\lambda_{\alpha/2}\sqrt{\hat{V}(\hat{\theta})}}{\hat{\theta}}$ – alternatiivne suhteline viga
- (5) $MSE(\hat{\theta}) = E(\hat{\theta} - \theta)^2 = V(\hat{\theta}) + B^2$, kus $B = E(\hat{\theta}) - \theta$ – keskmine ruutviga.
- (5') $M\hat{S}E(\hat{\theta})$ – keskmise ruutvea hinnang.

Valikuviga ja seega ka hinnangu täpsus sõltub:

- valikudisainist (sealhulgas valimimahust, $n = \sum_U I_i$)
- hinnangufunktsioonist

Valimimahu suurendamine suurendab hinnangu täpsust!

3.5.1 Valimimahu määramine

Valimimaht n määratakse vastavalt tellija poolt nõutavale hinnangu täpsusele.

Näiteks,

(1) hinnangu suhteline viga ei tohi ületada 2%, st

$$\frac{\sqrt{\hat{V}(\hat{t}_y)}}{\hat{t}_y} \leq 0.02 \text{ ehk } \hat{V}(\hat{t}_y) \leq (0.02\hat{t}_y)^2.$$

Nendest võrratustest saab määrata valimimahu, kui teame $\hat{V}(\hat{t}_y)$ ja \hat{t}_y hinnangulisi väärtusi, näiteks eelnevast uuringust või taustuuringust, ja vastavaid valemeid sõltuvalt valimimahust n .

(2) ILO (*International Labour Organization*) nõuab tööjõu uuringute läbi viimisel, et kasutatav disain oleks selline, et osakogumites, mis moodustavad 5% ÜK-st, ei ületaks hinnangu standardviga 6% hinnangust. Teiste sõnadega, osakogumites, mille maht $N_d = 0.05N$ on nõue hinnangutele järgmine:

$$\sqrt{\hat{V}(\hat{t}_d)} \leq 0.06\hat{t}_d.$$

Selline täpsus nõuab väga suurt valimimahtu.

3.6 Disainiefekt

- Aitab võrrelda disaine hinnangute täpsuse seisukohalt.
- Suur praktiline väärtus komplitseeritud disainide korral, mil pole võimalik leida hinnangute dispersioonivalemeid. Sel juhul disainiefekti ligikaudne teadmine (nt. eelnevate uuringute kogemustest) võimaldab sellistes olukordades hinnata dispersioone.

Def. Valikudisaini $p(s)$ disainiefekt on suhe

$$Def f_{p(s)}(\hat{t}_y) = \frac{V_{p(s)}(\hat{t}_y)}{V_{LJV}(\hat{t}_y)}.$$

LJV TTA on võetud disainiks, mille suhtes võrreldakse teisi disaine, kuna see on teoreetiliselt hästi läbitöötatud ja praktikas sageli kasutatav disain. Disainiefekt sõltub uuritavast tunnusest y_i ja valikudisainist $p(s)$

4 Hindamine lihtsa juhuvaliku korral, TTA

Olgu $I = (I_1, \dots, I_N)$ disaini vektor lihtsa juhuvaliku korral ning valimi maht olgu n . Sellisel juhul:

$$I \sim p(k) = \begin{cases} (C_N^n)^{-1}, & \text{kui } |k| = n; \\ 0, & \text{vastasel juhul.} \end{cases}$$

Tänu oma lihtsusele, LJV ja ka hinnangud selle disaini korral on väga hästi uuritud.

Kasutusvaldkond:

- Sageli on LJV osa mingist keerulisemast disainist (nt. kaheastmeline disain, kus 1.-l astmel valitakse vastavalt LJV-le suurimad objektid ehk klastrid (majad, tänavad, osariigid jne). Ja teisel astmel igas klastris rakendatakse oma (kõige sobivam) disain.
- Valemid, mis on välja töötatud LJV jaoks võivad olla rakendatud lähendina teiste disainide korral, näiteks hinnangu dispersiooni valem süstemaatilise valiku korral.

4.1 LJV disainikarakteristikud

Kõigepealt, esimest järku kaasamistõenäosus: $\pi_i = Pr(I_i = 1) = ?$ Sündmus $I_i = 1$ toimub kui vektor $I = (I_1, \dots, I_i, \dots, I_N)$ saab realisatsiooniks $k_i = 1$, $k = (k_1, \dots, 1, \dots, k_N)$. Kuna soovime, et valimimaht oleks n , siis selliseid võimalike realisatsioonide on C_{N-1}^{n-1} . Seega, saame

$$\pi_i = Pr(I_i = 1) = \sum_{k, k_i=1, |k|=n} p(k) = C_{N-1}^{n-1} (C_N^n)^{-1} = \frac{n}{N}.$$

Suhet n/N nimetatakse sageli valikusuhteks (*sampling fraction*) ja tähistatakse f -ga.

Analoogiliselt saame avaldise 2.-st järku kaasamistõenäosusele:

$$\pi_{ij} = Pr(I_i = 1, I_j = 1) = \sum_{k, |k|=1, k_i=k_j=1} p(k) = C_{N-2}^{n-2} (C_N^n)^{-1} = \frac{n(n-1)}{N(N-1)}.$$

Samuti läheb edaspidi vaja valikuindikaatorite dispersiooni ja kovariatsiooni:

$$\begin{aligned} \Delta_{ii} &= V(I_i) \stackrel{def}{=} E(I_i^2) - (EI_i)^2 = \dots \\ E(I_i^2) &= 1 \cdot Pr(I_i^2 = 1) + 0 \cdot Pr(I_i^2 = 0) = Pr(I_i = 1) = \pi_i \\ \dots &= \pi_i - \pi_i^2 = \pi_i(1 - \pi_i) = \frac{n}{N} \left(1 - \frac{n}{N}\right) \\ \Delta_{ij} &= Cov(I_i, I_j) \stackrel{def}{=} \underbrace{E(I_i I_j)}_{\pi_{ij}} - E(I_i)E(I_j) = \pi_{ij} - \pi_i \pi_j = \\ &= \frac{n(n-1)}{N(N-1)} - \left(\frac{n}{N}\right)^2 = \frac{n}{N} \left(\frac{n-1}{N-1} - \frac{n}{N}\right) = \\ &= \frac{n}{N} \frac{Nn - N - Nn + n}{N(N-1)} = -\frac{n}{N} \frac{N-n}{N} \frac{1}{N-1}. \end{aligned}$$

Kokkuvõttes saame LJV disainikarakteristikuid koondada järgmisesse tabe-

lisse:

$$\begin{aligned}
 f &= \frac{n}{N} && \text{valikusuhe;} \\
 \pi_i &= f && \text{esimest järku kaasamistõenäosus;} \\
 \pi_{ij} &= f \frac{n-1}{N-1}, i \neq j && \text{teist järku kaasamistõenäosus;} \\
 \Delta_{ii} &= f(1-f) && I_i \text{ dispersioon;} \\
 \Delta_{ij} &= -f(1-f) \frac{1}{N-1} && I_i, I_j \text{ kovariatsioon.}
 \end{aligned}$$

4.2 Hinnang kogusummale LJV korral

Meid huvitab parameeter $t = \sum_U y_i$. Üldisest hindamisteoreemist saame:

$$\hat{t} = \sum_U \frac{I_i y_i}{EI_i} = \sum_U \frac{I_i y_i}{\pi_i} = \frac{N}{n} \sum_U I_i y_i,$$

millele vastab järgmine tavapärase kuju:

$$\hat{t} = \frac{N}{n} \sum_s y_i = N\bar{y}.$$

Sellel hinnangul on olemas kaks tõlgendust:

$$\hat{t} = \begin{cases} N\bar{y} & \text{– valimikeskmine esindab kõiki väärtuseid ÜK-st;} \\ \sum_s w_i y_i & \text{– iga valimiäärtus } y_i \text{ esindab } w_i = N/n \text{ väärtust ÜK-st.} \end{cases}$$

Järgmisena võiksime leida saadud hinnangu dispersiooni ja dispersiooni hinnangu kasutades ÜHT. Kuid topeltsummad dispersiooni avaldises võivad osutada aeganõudvaks suurte andmestikke korral.

Kuna LJV on fikseeritud maguga disain, siis saame kasutada alternatiivset valemit, mida leidsime punktis 3.1.3 (teoreem 5).

$$\begin{aligned}
 V(\hat{t})_{alt} &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \Delta_{ij} \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 \underbrace{\Delta_{ij} = -f(1-f) \frac{1}{N-1}} \\
 &= \frac{1}{2} f(1-f) \frac{1}{f^2} \frac{1}{N-1} \sum_i \sum_j (y_i - y_j)^2.
 \end{aligned}$$

Lisades $\pm\bar{Y}$ ja avaldades ruutu, paneme tähele, et

$$\sum_i \sum_j (y_i - \bar{Y})^2 = N \sum_{i=1}^N (y_i - \bar{Y})^2,$$

$$\sum_{i=1}^N (y_i - \bar{Y}) = \sum_{i=1}^N y_i - \underbrace{N\bar{Y}}_{\sum_{i=1}^N y_i} = 0.$$

Järelikult,

$$\sum_i \sum_j (y_i - y_j)^2 = 2N \sum_{i=1}^N (y_i - \bar{Y})^2.$$

Lõplikult saame,

$$V(\hat{t}) = \frac{1}{2}(1-f) \frac{1}{f} 2N \underbrace{\frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{Y})^2}_{S_y^2} = N^2(1-f) \frac{S_y^2}{n},$$

kus S_y^2 on tunnuse y üldkogumi dispersioon. NB! Pole enam topitsummasid!

Analoogiliselt dispresiooniga saame lihtsustada ka dispersiooni hinnangut:

$$\hat{V}(\hat{t})_{alt} = -\frac{1}{2} \sum_{i,j \in s} \underbrace{\frac{\Delta_{ij}}{\pi_{ij}}}_{-\frac{1-f}{n-1}} \left(\frac{y_i}{\pi_i} - \frac{y_j}{\pi_j} \right)^2 = \frac{1-f}{2(n-1)f^2} \sum_{i,j \in s} (y_i - y_j)^2.$$

Nüüd saame kasutada, et

$$\sum_{i,j \in s} (y_i - \bar{y})^2 = n \sum_{i \in s} (y_i - \bar{y})^2,$$

$$\sum_{i \in s} (y_i - \bar{y}) = \sum_s y_i - n\bar{y} = 0.$$

Ja lõplikult saame:

$$\hat{V}(\hat{t}) = \frac{1-f}{2(n-1)f^2} 2n \sum_s (y_i - \bar{y})^2 = \quad (29)$$

$$= N^2(1-f) \frac{1}{n} \underbrace{\frac{1}{n-1} \sum_s (y_i - \bar{y})^2}_{s_y^2} = N^2(1-f) \frac{s_y^2}{n}, \quad (30)$$

kus s_y^2 on tunnuse y valimi dispersioon.

Teoreem 6. Lihtsa juhuvaliku TTA korral nihketa hinnang \hat{t} summale $t = \sum_U y_i$ avaldub järgmiselt:

$$\hat{t} = \frac{N}{n} \sum_U I_i y_i = \frac{N}{n} \sum_s y_i,$$

ehk alternatiivselt

$$\hat{t} = N\bar{y}.$$

Hinnangu dispersioon on järgmine:

$$V(\hat{t}) = N^2(1-f) \frac{S_y^2}{n}$$

ja dispersiooni hinnang:

$$\hat{V}(\hat{t}) = N^2(1-f) \frac{s_y^2}{n},$$

kus

$$\begin{aligned} f &= \frac{n}{N} \text{ on valikusuhe,} \\ \bar{y} &= \frac{1}{n} \sum_s y_i \text{ valimikeskmine,} \\ S_y^2 &= \frac{1}{N-1} \sum_U (y_i - \bar{Y})^2 \text{ tunnuse } y \text{ ÜK dispersioon,} \\ s_y^2 &= \frac{1}{n-1} \sum_s (y_i - \bar{y})^2 \text{ tunnuse } y \text{ valimidispersioon.} \end{aligned}$$

Järeldus 6.1 Hinnangud \hat{Y} keskmisele $\bar{Y} = t/N$ on järgmised:

$$\begin{aligned}\hat{Y} &= \bar{y}, \\ V(\hat{Y}) &= (1-f)S_y^2/n, \\ \hat{V}(\hat{Y}) &= (1-f)s_y^2/n.\end{aligned}$$

Võrdle klassikalise statistika valemitega!

Klassikalises statistikas on y_i sõltumatud sama jaotusega valimis ning

$$V(\hat{Y}) = \frac{S_y^2}{n}, \text{ ja } \hat{V}(\hat{Y}) = \frac{s_y^2}{n}.$$

LJV TTA korral kui me eemaldame ühe y_i ÜK-st, siis saame hoopis teise ÜK jaotuse. Järgmise elemendi valik toimub juba teise jaotuse järgi. Valikud on omavahel negatiivselt korreleeritud:

$$\Delta_{ij} = -f(1-f)\frac{1}{N-1}.$$

Seepärast ka keskmise hiinangu dispersioonil olemas nn "lõpliku ÜK-i korrigeerimiskordaja" $(1-f)$.

Kui valikusuhe on väike, siis ka valemid on ligikaudselt võrdsed klassikalise statistika omadega.

Jäta meelde! Klassika statistika tarkvara pakettid töötavad eeldusel, et y_i on sõltumatute samast jaotuseset. Seepärast neid ei saa otseselt kasutada valikuuringute andmetele. Mida keerulisem on disain, seda rohkem erinevad hinnangute dispersioonid.

SUDAAN	}	Spetsiaalne tarkvara valikuuringute teostamiseks.
WESVAR		
CLAN		
SAS 8.1		
R, pakett SAMPLING		

4.3 Kovariatsioon kahe hinnangu vahel LJV TTA korral

Olgu meil kaks hinnangut $\hat{t}_y = N\bar{y}$ ja $\hat{t}_x = N\bar{x}$. Üldiselt TTA disainide jaoks avaldub kovariatsioon kahe hinnangu vahel järgmiselt (vt punkt 3.3):

$$Cov(\hat{t}_y, \hat{t}_x) = \sum_i \sum_j \Delta_{ij} \frac{y_i x_j}{\pi_i \pi_j} = \sum_i \Delta_{ii} \frac{y_i x_i}{\pi_i^2} + \sum_{i \neq j} \sum_j \frac{y_i x_j}{\pi_i \pi_j}.$$

LJV TTA korral:

$$\begin{aligned} Cov(\hat{t}_y, \hat{t}_x) &= \frac{N^2}{n^2} \left[f(1-f) \sum_i y_i x_i - f(1-f) \frac{1}{N-1} \sum_{i \neq j} \sum_j y_i x_j \right] = \\ &= \frac{N}{n} (1-f) \frac{1}{N-1} \left[(N-1) \sum_i x_i y_i - \underbrace{\sum_{i \neq j} \sum_j y_i x_j}_{-[(\sum_i y_i)(\sum_i x_i) - \sum_i y_i x_i]} \right] = \\ &= \frac{N}{n} (1-f) \frac{1}{N-1} \left[N \sum_i y_i x_i - t_y t_x \right] = N^2 (1-f) \frac{1}{N-1} \underbrace{\left[\sum_i y_i x_i - N\bar{Y}\bar{X} \right]}_{S_{yx}} / n. \end{aligned}$$

Lõplikult saame:

$$Cov(\hat{t}_y, \hat{t}_x) = N^2 Cov(\hat{Y}, \hat{X}) = N^2 (1-f) S_{yx} / n,$$

kus

$$S_{yx} = \frac{1}{N-1} \sum_U (y_i - \bar{Y})(x_i - \bar{X})$$

on tunnuste y ja x kovariatsioon ÜK-s.

Punktist 3.3 saab ka tuletada hinnangu kovariatsioonile:

$$\hat{C}ov(\hat{t}_y, \hat{t}_x) = \sum_{i,j \in s} \sum_{\pi_{ij}} \frac{\Delta_{ij} y_i x_i}{\pi_i \pi_i} \stackrel{LJVTTA}{=} N^2 (1-f) \frac{S_{yx}}{n},$$

kus

$$s_{yx} = \frac{1}{n-1} \sum_s (y_i - \bar{y})(x_i - \bar{x})$$

on valimi kovariatsioon y ja x vahel.

Ülesanne Tuletada $Cor(\hat{t}_y, \hat{t}_x)$. Kommenteerida!

4.4 Suhtehinnang LJV TTA korral

Lihtsa juhuvaliku korral avaldub kahe kogusumma suhte $R = \frac{t_y}{t_x}$ hinnang kujul

$$\hat{R} = \frac{\hat{t}_y}{\hat{t}_x} = \frac{\bar{y}}{\bar{x}} = \frac{\sum_s y_i}{\sum_s x_i}.$$

Hinnangu \hat{R} ligikaudne dispersioon avaldub üldjuhul järgmiselt:

$$AV(\hat{R}) = \frac{1}{t_x^2} [V(\hat{t}_y) + R^2 V(\hat{t}_x) - 2RCov(\hat{t}_y, \hat{t}_x)].$$

Sellest valemist saame LJV TTA korral

$$AV(\hat{R}) = \frac{1-f}{n\bar{X}^2} (S_y^2 + R^2 S_x^2 - 2RS_{xy}),$$

kus S_y^2, S_x^2 on ÜK dispersioonid (vt teoreemi 6) ja S_{xy} on üldkogumi kovariatsioon vaadeldud punktis 4.3.

Suhtehinnangu dispersiooni hinnang avaldub lihtsa juhuvaliku korral järgmiselt:

$$\hat{V}(\hat{R}) = \frac{1-f}{n\bar{x}^2} (s_y^2 + \hat{R}^2 s_x^2 - 2\hat{R}s_{yx}), \quad (31)$$

kus s_y^2, s_x^2 on valimi dispersioonid ja s_{yx} on valimi kovariatsioon.

Praktilises töös, näiteks statistikaametites armastatakse suhte dispersiooni hindamisel kasutada teist lähenemist, mis seisneb sobivate uute tunnuste moodustamises. Eespool saime lähendava lineaarse valemi suhtele

$$\hat{R} \approx R + \frac{1}{t_x} \sum_U I_i u_i, \quad \text{kus } u_i = y_i - Rx_i.$$

Dispersioon tuleb teisest liidetavast,

$$V(\hat{R}) = \frac{1}{\hat{t}_x^2} V\left(\sum_U I_i u_i\right).$$

Lähtudes üldistest dispersioonivalemitest näitasime eespool, et LJV TTA korral $\hat{t}_y = \sum_U I_i y_i$ dispersioon ja dispersiooni nihketa hinnang on

$$V(\hat{t}_y) = N^2(1-f)S_y^2/n \text{ ja } \hat{V}(\hat{t}_y) = N^2(1-f)s_y^2/n.$$

Meie jaoks tulevad need valemid nüüd väljendada tunnuse u_i kaudu. Sõnas-tame teoreemina.

Teoreem 7. Suhte $\hat{R} = \frac{\hat{t}_y}{\hat{t}_x} = \frac{\sum_s y_i}{\sum_s x_i}$ ligikaudne dispersioon avaldub LJV TTA korral valemiga

$$V(\hat{R}) \approx \frac{N^2(1-f)}{\hat{t}_x^2} S_u^2/n,$$

ja dispersioonihinnang valemiga

$$\hat{V}(\hat{R}) = \frac{N^2(1-f)}{\hat{t}_x^2} s_u^2/n, \quad (32)$$

kus

$$S_u^2 = \sum_{i=1}^N (u_i - \bar{U})^2 / (N-1), \quad s_u^2 = \sum_s (u_i - \bar{u})^2 / (n-1).$$

Valemiga (32) on dispersioonihinnangut oluliselt lihtsam leida kui valemiga (31)

Ülesanne. Avaldada \bar{U} ja \bar{u} .

4.5 Hindamine osakogumites LJV TTA korral

Kui ÜK on mingi tunnuse väärtuste järgi jagatud osadeks ehk osakogumiteks, siis huvitavad meid nende osakogumite mahud - nii absoluutselt kui suhtelised, kogusummad, keskmised ja suhted osakogumites. Olgu üldkogumiks kõik õpilased ja osakogumiks 1. klassi õpilased. Meid huvitavad järgmised osakogumi näitajad:

Näiteks:

- osakogumi maht ehk kõigi 1. klassi õpilaste arv;
- matemaatika õppimisele kulutatud summaarne aeg;
- keskmine matemaatika õppimisele kuluv aeg õpilase kohta;
- matemaatika õppimisele kuluv aeg osakaaluna kodutöödele kuluvast ajast.

Vaatame kõigepealt osakogumi mahu ja osakaalu hindamist. Osakaalu esitatakse tavaliselt protsentides. Olgu $U_d \subset U$ meid huvitav osakogum, N_d, N on vastavad mahud.

Osakaal on $P_d = N_d/N$. Osavalimi s_d maht on n_d , osakaal valimis on $p_d = n_d/n$.

Iga üldkogumi objektiga seotakse osakogumi indikaator:

$$z_i = \begin{cases} 1, & i \in U_d \\ 0, & \text{vastasel juhul.} \end{cases}$$

Selle tunnuse kogusumma ja keskmine on osakogumi maht N_d ja suhteline maht P_d :

$$N_d = t_z = \sum_U z_i, \quad P_d = \frac{N_d}{N} = \bar{Z} = \frac{t_z}{N}.$$

Nüüd saame rakendada neile ÜHT LJV erijuhul:

$$\begin{aligned} \hat{N}_d &= \hat{t}_z = N\bar{z} = Np_d, \text{ kus } p_d \text{ on valimi } s_d \text{ osakaal;} \\ \hat{P}_d &= \frac{\hat{t}_z}{N} = \bar{z} = p_d. \end{aligned}$$

Hinnangu dispeersiooni saamiseks vaatleme esmalt S_z^2 :

$$S_z^2 = \frac{1}{N-1} \left(\sum_U \underbrace{z_i^2}_{z_i^2=z_i, z_i \in \{0,1\}} - N\bar{Z}^2 \right) = \frac{1}{N-1} \left(\sum_U \underbrace{z_i}_{NP_d} - NP_d^2 \right) = \frac{N}{N-1} P_d(1-P_d).$$

Järelikult hinnangute \hat{N}_d ja \hat{P}_d dispersioonid ja dispersiooni hinnangud avalduvad järgmiselt (näita!):

$$\begin{aligned} V(\hat{N}_d) &= \frac{N^3}{n}(1-f)\frac{1}{N-1}P_d(1-P_d), & V(\hat{P}_d) &= \frac{1-f}{n}\frac{N}{N-1}P_d(1-P_d) \\ \hat{V}(\hat{N}_d) &= N^2\frac{1-f}{n-1}p_d(1-p_d), & \hat{V}(\hat{P}_d) &= \frac{1-f}{n-1}p_d(1-p_d) \end{aligned}$$

Tunnuse y kogusumma osakogumis U_d on $t_d = \sum_{U_d} y_i$. Selleks, et oleks võimalik rakendada ÜHT hinnangute saamiseks, peame esitama t_d ÜK kogusumma kaudu. Selleks kasutame jällegi binaarset tunnust z :

$$t_d = \sum_{U_d} y_i = \sum_U y_i z_i = \sum_U y'_i, \text{ kus } y'_i = z_i y_i.$$

Teoreemist 6 saame summa hinnangule LJV TTA korral kuju:

$$\hat{t}_d = N\bar{y}' = \frac{N}{n} \sum_s y'_i = \frac{N}{n} \sum_{s_d} y_i.$$

Dispersioonide $V(\hat{t}_d)$ ja $\hat{V}(\hat{t}_d)$ leidmiseks asendame tunnuse y_i tunnusega y'_i Teoreemis 6.

Ülesanne. Kirjuta välja saadud hinnangu dispersioon ja dispersiooni hinnang.

Osakogumi keskväärtuse $\bar{Y}_d = t_d/N_d$ hindamiseks saab kasutada järeldust 6.1:

$$\hat{Y}_d = \frac{\hat{t}_d}{N_d} = \frac{N}{N_d} \frac{1}{n} \sum_{s_d} y_i.$$

Pane tähele, et saadud hinnang pole osakogumi valimikeskmine! Keskväärtuse hindamiseks on tavaliselt alternatiivne hinnang ehk suhtetüüpi hinang parem:

$$\hat{Y}_{d,alt} = \frac{\hat{t}_d}{\hat{N}_d}.$$

Vaatame, mis kuju see võtab LJV TTA korral. Paneme tähele, et $\hat{N}_d = N/n \sum_s z_i = (N \cdot n_d)/n$. Siit saame, et

$$\hat{Y}_{d,alt} = \frac{N/n \sum_{s_d} y_i}{N \cdot n_d/n} = \frac{1}{n_d} \sum_{s_d} y_i = \bar{y}_d.$$

Seega, suhtetüüpi keskmise hinnang osakogumis on valimi keskmine selles osakogumis. Kui osakogumi maht N_d on teada, siis on parem kogusumma hinnang osakogumis $N_d \hat{Y}_{d,alt}$.

Üldjuhul huvitab meid järgmine osakogumi suhe, mille aga saame esitada suhtena kogu üldkogumis uute tunnuste abil:

$$R_d = \frac{\sum_{U_d} y_i}{\sum_{U_d} x_i} = \frac{\sum_U y'_i}{\sum_U x'_i},$$

kus $y'_i = z_i y_i$ ja $x'_i = z_i x_i$. Hinnangu sellele suhtele saame lugeja ja nimetaja nihketa hindamise teel:

$$\hat{R}_d = \frac{N \bar{y}'}{N \bar{x}'} = \frac{\sum_s y'_i}{\sum_s x'_i} = \frac{\sum_{s_d} y_i}{\sum_{s_d} x_i}.$$

Dispersioonivalemid $V(\hat{R}_d)$ ja $\hat{V}(\hat{R}_d)$ järelduvad Teoreemist 7, milles tuleb kasutada osakogumitunnuseid y'_i ja x'_i .

5 Hindamine lihtsa juhuvaliku TGA korral

LJV TGA valiku korral tehakse ÜK-s U n valikut tõenäosusega

$$p_i = \frac{1}{N}, \quad \sum_{i=1}^N p_i = 1.$$

Iga kord kui objekt on valitud, pannakse see üldkogumisse tagasi.

LJV TGA korral on valikuindikaatorid binoomjaotusega: $I_i \sim \text{Bin}(n, \frac{1}{N})$.

Disaini vektori I jaotuseks on multinomiaalne jaotus:

$$p(k) = \Pr(I = k) = \begin{cases} \frac{n!}{N^n \prod_{i=1}^N k_i!}, & \text{kui } |k| = n, \\ 0, & \text{vastasel juhul.} \end{cases}$$

Binoom- ja multinoomjaotuse karakteristikud on hästi tuntud:

$$\begin{aligned} E(I_i) &= np_i = \frac{n}{N} \\ V(I_i) &= np_i(1 - p_i) = \frac{n}{N}\left(1 - \frac{1}{N}\right), \\ Cov(I_i, I_j) &= -np_i p_j = -\frac{n}{N^2}. \end{aligned}$$

Kasutusvaldkonnad:

- valikudisainina ei kasutata;
- valemid, tuletatud LJV TGA jaoks on tavaliselt lihtsa ja ilusa kujuga, neid saab kasutada sageli lähendina teiste disainide juures sobivas olukorras;
- LJV TGA on tähtis nn "taasvaliku"teoorias, kus saadud valimist võetakse korduvalt omakorda valimid kasutades LJV TGA ja selle protseduuri abil leitakse hinnangu dispersiooni hinnang.

Kogusumma hinnangu ja selle dispersiooni tuletamiseks kasutame teoreemi 4 punktist 3.1.2(hinnangud multinomiaaldisaini korral üldjuhul). Dispersiooni hinnangu saamiseks aga kasutame aga alternatiivset valemit, kuna LJV on fikseeritud mahuga disainid ja Sen-Yates-Grundy hinnang on stabiilsem (varieeruvuse mõttes, samuti ei võta ta negatiivseid väärtuseid).

Olgu $\hat{t} = \sum_U \frac{I_i y_i}{E(I_i)}$ on kogusumma $t = \sum_U y_i$ hinnang. Multinomiaaldisaini korral alternatiivne dispersioonihinnang avaldub järgmiselt:

$$\begin{aligned} \hat{V}(\hat{t}) &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N I_i I_j \frac{\Delta_{ij}}{E(I_i I_j)} \left(\frac{y_i}{E I_i} - \frac{y_j}{E I_j} \right)^2 = \\ &= -\frac{1}{2} \sum_i \sum_j I_i I_j \frac{-np_i p_j}{n(n-1)p_i p_j} \left(\frac{y_i}{np_i} - \frac{y_j}{np_j} \right)^2 \Rightarrow \\ \hat{V}(\hat{t}) &= \frac{1}{2} \frac{1}{n-1} \sum_i \sum_j I_i I_j \left(\frac{y_i}{np_i} - \frac{y_j}{np_j} \right)^2. \end{aligned} \quad (33)$$

Avaldame ruutu ja lihtsustame:

$$\begin{aligned} \sum_i \sum_j I_i I_j \frac{y_i^2}{n^2 p_i^2} &= \sum_i I_i \frac{y_i^2}{n^2 p_i^2} \underbrace{\sum_{j=1}^2 I_j}_n = \frac{1}{n} \sum_i I_i \frac{y_i^2}{p_i^2}, \\ \sum_i \sum_j I_i I_j \frac{y_i}{n p_i} \frac{y_j}{n p_j} &= \underbrace{\sum_i I_i \frac{y_i}{n p_i}}_{\hat{t}} \underbrace{\sum_j I_j \frac{y_j}{n p_j}}_{\hat{t}} = \hat{t}^2 \end{aligned}$$

Lõpuks saame,

$$\begin{aligned} \hat{V}(\hat{t}) &= \frac{1}{2n-1} \frac{1}{n} \left[\frac{2}{n} \sum_i I_i \frac{y_i^2}{p_i^2} - 2\hat{t}^2 \right] = \\ &= \frac{1}{n(n-1)} \left[\sum_{i=1}^N I_i \frac{y_i^2}{p_i^2} - n\hat{t}^2 \right], \end{aligned}$$

vastava hinnanguga

$$\hat{V}(\hat{t}) = \frac{1}{n(n-1)} \left[\sum_s k_i \frac{y_i^2}{p_i^2} - n\hat{t}^2 \right].$$

Teoreem 7. Lihtsa juhusliku valiku tagasipanekuga korral nihketa hinnang ÜK kogusummale $t = \sum_U y_i$ avaldub järgmiselt:

$$\hat{t} = \frac{N}{n} \sum_U I_i y_i,$$

vastava punktihinnanguga

$$\hat{t} = N\bar{y}.$$

Hinnangu \hat{t} dispersioon on järgmine:

$$V(\hat{t}) = \frac{N(N-1)}{n} S_y^2,$$

ja dispersiooni hinnangufunktsioon:

$$\hat{V}(\hat{t}) = \frac{N^2}{n(n-1)} \left[\sum_{i_1}^N I_i y_i^2 - n\bar{y}^2 \right].$$

Viimasele avaldisele vastav punktihinnang on järgmine:

$$\hat{V}(\hat{t}) = \frac{N^2}{n} s_y^2,$$

kus

$$\begin{aligned} \bar{y} &= \frac{1}{n} \sum_s y_i, \\ \bar{Y} &= \frac{1}{N} \sum_U y_i, \\ S_y^2 &= \frac{1}{N-1} \sum_U (y_i - \bar{Y})^2, \\ s_y^2 &= \frac{1}{n-1} \sum_s (y_i - \bar{y})^2. \end{aligned}$$

Hinnangud keskväärtusele LJV TGA korral avalduvad järgmiselt:

$$\begin{aligned} \hat{Y} &= \bar{y}, \\ \hat{V}(\hat{Y}) &= \frac{s_y^2}{n}. \end{aligned}$$

Saime klassikalise statistika tulemusi!

Võrdleme SI ja SIR (LJV TTA ja LJV TGA) omavahel:

$$V_{SI}(\hat{t}) = N^2(1-f) \frac{S_y^2}{n}, \text{ kus } f = \frac{n}{N} \text{ ja } V_{SIR}(\hat{t}) = N(N-1) \frac{S_y^2}{n}.$$

Juhul kui $n = N$, $V_{SI}(\hat{t}) = 0$, kuid $V_{SIR}(\hat{t}) \neq 0$!

Üldiselt, $V_{SIR}(\hat{t}) \geq V_{SI}(\hat{t})$. Võrdusmärk kehtib kui $n = 1$ ja $n = N - 1$. Lihtne juhuvalik tagasipanekuga on vähem efektiivne kui tagasipanekuta lihtne juhuvalik.

Disainiefekt:

$$Def_{SIR} = \frac{V_{SIR}(\hat{t})}{V_{SI}(\hat{t})} = \frac{N(N-1)S_y^2/n}{N^2(1-f)S_y^2/n} \approx \frac{1}{1-f}.$$

Mida suurem on valikusuhe f (st mida lähedasem on ta 1-le), seda vähem on SIR efektiivne.

5.1 Isekaaluvad disainid

Isekaaluvate disainide korral kehtib:

$$E(I_i) \equiv const.$$

Kui lisaks disain on fikseeritud mahuga n , siis $\sum_{i=1}^N E(I_i) = n$ ja siit järeldub, et

$$E(I_i) \equiv \frac{n}{N}.$$

Sellisel juhul nihketa hinnang \hat{t} kogusummale on

$$\hat{t} = \sum_U \frac{I_i y_i}{E(I_i)} = \frac{N}{n} \sum_U I_i y_i,$$

mis tähendab, et hinnang põhineb valimikeskmisel:

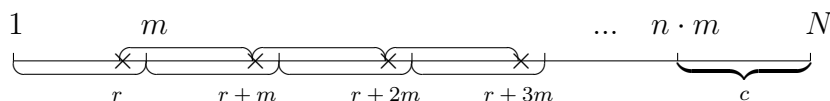
$$\hat{t} = N\bar{y}.$$

Ja see omakorda tähendab, et \hat{t} keskmiist hindab valimikeskmine, \hat{t} osakaalu - valimi osakaal, ... \Leftrightarrow valimikarakteristikud esindavad \hat{t} karakteristiku. See on aga küllaltki mugav ning võimaldab kasutada hindamisülesannetes tarkvara, mis on mõeldud klassikalise statistika jaoks.

Tähelepanu! Kogusumma hinnangu disperioonid on siiski üldjuhul erinevad eri disainide korral!

6 Süstemaatiline valik

Olgu $U = 1, \dots, N$.



Süstemaatilise valiku korral võetakse esimene element valimisse juhuslikult m esimese elemendi hulgast (võrdse tõenäosusega). Valimit moodustavad see esimene element pluss iga m -s element freimist.

Kokku on võimalik saada m erinevat valimit. Iga sellise valimi saamise tõenäosus on $1/m$.

Olgu $I = (i_1, \dots, i_N)$ valikuvektor freimis U . Süstemaatilise valiku korral on sellise vektori realisatsiooniks k vektor elementidega 0 ja 1, kus 1 esineb iga m sammu tagant. Seega, on vektori I jaotus järgmine:

$$p(k) = \Pr(I = k) = \begin{cases} 1/m, & \text{kui 1 ilmub esimese } m \text{ hulgas;} \\ 0, & \text{vastasel juhul.} \end{cases}$$

Süstemaatiline disain on TTA disain. Kasutusvaldkonnad:

1. *SÜ on lihtsalt teostatav jooksval valikul ja seetõttu ta on vähem tundlik intervjuerijate subjektiivsete vigade suhtes kui LJV või KV (kihtvalik). Eriti kui korralik freim ei ole kättesaadav.*

Näiteks, ostjate lihtsa juhusliku valimi mahuga $n = 50$ korjamine tänavanurgal on üsna keerukas. Intervjuerija ei saa otsustada, milliseid ostjaid võtta valimisse, sest ÜK maht N ei ole teada kuni kõik ostjad on ümber nurga keeranud. Seevastu intervjuerija võib kasutada SÜ-t ja võtta valimisse näiteks iga 20. ostja kuni nõutava mahuga valim on saadud. See protseduur on lihtne isegi kogenematu intervjuerija jaoks.

Need küsitlajad, kes küsitlevad inimesi liikumisel, kasutavadki väga sageli SÜ-t. Nad võivad küsitleda iga 20-nda inimese kassa juures toidu maitse või värvuse kohta. Iga 10-s isik, kes siseneb bussi võib olla küsitatud bussiteeninduse kohta. Samuti metsavahid võivad võtta maatükkide süstemaatilist

valimit, ning süstemaatiliselt valida puud ise, et uurida haigete puude osakaalu. Seetõttu on SÜ väga populaarne valikudisain.

2. SÜ võib anda täpsema informatsiooni kui LJV sama maksmuse korral.

Süstemaatiline valim on reeglina "ühtlasem" ja seega annab rohkem informatsiooni ÜK kohta kui sama mahuga lihtne juhuslik valim. Näiteks, me tahame võtta SÜ valimit mahuga $n = 200$ panga maksekviitungite ÜK-st mahuga $N = 1000$, selleks, et hinnata korrektselt täidetud kviitungite osakaalu. Selleks võtame juhuslikult ühe kviitungi 5-st esimesest määramaks alguspunkti (näiteks, number 3), ja seejärel võtame iga viienda kviitungi.

Oletame, et suurem osa esimesest 500st kviitungist olid täidetud korrektselt, järgmised 500 aga olid kõik täidetud valesti (näiteks, pangateenindaja kogenematus tõttu). LJV korral võib valimisse ($n = 200$) sattuda liiga palju (võimalik, et kõik) kviitunge esimesest (või teisest) osast kviitungitest. See annab aga ebatäpse hinnangu osakaalule. Seevastu SÜ valib võrdse kviitungite arvu mõlemast gruppist ja annab parema hinnangu valesti täidetud kviitungite osakaalule.

SÜ korral igal objektil on olemas võimalus sattuda valimisse, st $\pi_i = \Pr(I_i = 1) > 0$. Disaini puudus on aga see, et mõned objektid ei saa korruga sattuda valimisse, st mõned 2. järku kaasamistõenäosused on võrdsed 0-ga, $\pi_{ij} = 0$. See aga omakorda tähendab, et pole võimalik leida hinnangute teoreetilist dispersiooni.

Valimimaht, \mathbf{n} , on SÜ korral juhuslik ja on määratud sammuga m . SÜ korral kehtib:

$$N = nm + c, \quad 0 \leq c < m.$$

Seega, realiseerunud valimimaht, tähistame n_s on:

$$n_s = \begin{cases} n + 1, & \text{kui } r \leq c; \\ n, & \text{kui } r > c. \end{cases}$$

Leiame ka 1. ja 2. järku kaasamistõenäosused:

$$\pi_i = \Pr(I_i = 1) = \sum_{k, k_i=1} p(k) = \frac{1}{m},$$

kuna on võimalik ainult üks selline valim k , mis sisaldaks i ndas positsioonis 1;

$$\pi_{ij} = \Pr(I_i = 1, I_j = 1) = \sum_{k, k_i=1, k_j=1} p(k) = \begin{cases} 1/m, & \text{kui vahe } i \text{ ja } j \text{ vahel on sammu } m \text{ kordne;} \\ 0, & \text{vastasel juhul.} \end{cases}$$

Kuna valimimahud väga ei varieeru, siis on huvitav leida keskmist valimi-mahtu:

$$E(\mathbf{n}) = E \left[\sum_{i=1}^N I_i \right] = \sum_{i=1}^N E(I_i) = \sum_{i=1}^N \pi_i = \frac{N}{m} = \frac{nm + c}{m} = n + \frac{c}{m}.$$

6.1 Hindamine SÜ korral

Kogusumma $t = \sum_U y_i$ nihketa hinnang ÜHT-i järgi on järgmine:

$$\hat{t} = \sum_U \frac{I_i y_i}{\pi_i} = m \sum_U I_i y_i,$$

millele vastab järgmine punktihinnang:

$$\hat{t} = m \sum_s y_i.$$

Kuna on võimalik saada kokku m erinevat valimit (olenevalt alguspunktist r), siis on ka võimalik saada kokku m erinevat hinnangut ÜK kogusummale, tähistame $\hat{t}_1, \dots, \hat{t}_m$:

$$\hat{t}_r = m \sum_{s_r} y_i, \quad r = 1, \dots, m.$$

Nüüd saame kirja panna $V(\hat{t})$ ilma ÜHT-ta, kasutades diskreetse juhusliku suuruse dispersiooni definitsiooni:

$$V(\hat{t}) = \sum_{i=1}^m (\hat{t}_r - \underbrace{E\hat{t}}_t)^2 \underbrace{\Pr(\hat{t} = \hat{t}_r)}_{1/m} = \frac{1}{m} \sum_{i=1}^m (\hat{t}_r - t)^2.$$

Näeme, et $V(\hat{t})$ sõltub sellest, kuidas varieeruvad \hat{t}_r kogusumma t ümber.

Dispersioon $V(\hat{t})$ on teoreetiline, seda ei saa välja arvutada, kuna praktikas on meil olemas ainult üks valim s_r ja ainult üks hinnang \hat{t}_r . Kuid me saame seda teoreetilist hinnangut kasutada SÜ uurimiseks.

Kõigepealt, lihtsustame situatsiooni ja eeldame, et $c = 0$, st $N = nm$.

Sel juhul $\hat{t}_r = m \sum_{s_r} y_i = \frac{N}{n} \sum_{s_r} y_i = N\bar{y}_r$. Kuna $t = N\bar{Y}$, me saame kirjutada teoreetilise dispersiooni $V(\hat{t})$ ümber järgmiselt:

$$V(\hat{t}) = Nn \sum_{r=1}^m (\bar{y}_r - \bar{Y})^2.$$

Seega, varieeruvus $V(\hat{t})$ sõltub valimikeskmiste varieeruvusest. Me soovime, et see varieeruvus oleks väike, see tagaks väikese dispersiooni $V(\hat{t})$. Kuna teisi valimeid pole, siis ka pole võimalik midagi otsustada selle varieeruvuse kohta. Antud olukorras saame kasutada ANOVA lahutust (tunnuse koguvarieeruvus grupisisese ja gruppidevahelise varieeruvuse kaudu):

$$\begin{aligned} SST &= \sum_U (y_i - \bar{Y})^2 = \sum_{r=1}^m \sum_{i \in s_r} (y_i - \bar{y}_r + \bar{y}_r - \bar{Y})^2 = \\ &= \underbrace{\sum_{r=1}^m \sum_{i \in s_r} (y_i - \bar{y}_r)^2}_{SSW} + \underbrace{n \sum_{r=1}^m (\bar{y}_r - \bar{Y})^2}_{SSB} = SSW + \frac{1}{N} V(\hat{t}). \end{aligned}$$

Fikseeritud ÜK korral on uuritava tunnuse varieeruvus, SST (*Sum of squares total*), samuti fikseeritud. Selleks, et saada väiksema $V(\hat{t})$, SSW (*Sum of squares within groups*) peab olema võimalikult suur. Ja see omakorda tähendab, et tunnus y varieeruvus valimis s_r peab olema suur \Rightarrow tunnus peab olema valimis s_r võimalikult heterogeenne. Järelikult, dispersiooni $V(\hat{t})$ suurus sõltub objektide järjestusest loendis.

Hea järjestus on järgmine:

- y väärtused, mis asuvad üksteisest kaugusel m peavad olema võimalikult erinevad;
- seda saab saavutada järjestades freimi väärtuseid kas uuritava tunnuse või sellega korreleeruva tunnuse väärtuste järgi.

Halb järjestus:

- loendis esineb väärtuste tsüklilise perioodiga m ; sellisel juhul tunnuse varieeruvus valimis on väike.

SÜ korral pole võimalik saada nihketa hinnangut kogusumma hinnangu dispersioonile, $\hat{V}(\hat{t})$. Sel juhul saab kasutada mõnda nihkega hinnangut, näiteks SI hinnangut:

$$\hat{V}(\hat{t}) = N^2(1-f) \frac{s_y^2}{n}.$$

Juhul, kui ÜK on halvasti järjestatud, siis s_y^2 võib osutuda liiga väikeseks ja sel juhul $\hat{V}(\hat{t})$ võib tegeliku dispersiooni alahinnata.

6.2 SÜ disaini efekt

Eelmises punktis näitsime, et uuritava tunnuse koguvarieeruvuse ÜK-s on võimalik esitada järgmiselt:

$$SST = SSW + SSB = SSW + \frac{1}{N}V(\hat{t}),$$

kus $SSB = \text{Sum Square Between}$ on varieeruvus gruppide vahel.

$$\begin{aligned} \Rightarrow V_{SY}(\hat{t}) &= N(SST - SSW) = N \cdot SST \left(1 - \frac{SSW}{SST}\right) = \\ &= N(N-1)S_y^2 \left(1 - \frac{SSW}{SST}\right). \end{aligned}$$

SÜ disainiefekt:

$$\begin{aligned} def f(SY) &= \frac{V_{SY}(\hat{t})}{V_{SI}(\hat{t})} = \frac{N(N-1)S_y^2 \left(1 - \frac{SSW}{SST}\right)}{N^2(1-f) \frac{S_y^2}{n}} = \\ &= \frac{(N-1)n}{N(1-f)} \left(1 - \frac{SSW}{SST}\right). \end{aligned}$$

Süsteemaatiline valik on efektiivsem kui lihtne juhuvalik siis, kui $def f(SY) < 1$. See aga tähendab järgmist võrratust:

$$\frac{(N-1)n}{N(1-f)} \left(1 - \frac{SSW}{SST}\right) < 1$$

$$\left(1 - \frac{SSW}{SST}\right) < \frac{N(1-f)}{(N-1)n}$$

$$\frac{SSW}{SST} > \frac{N(n-1)}{(N-1)n}$$

Tähistame

$$S_w^2 = \frac{SSW}{N-m}$$

- valimite sisene hajuvus. Arvestades, et

$$N-m = N - \frac{N}{m} = \frac{N(n-1)}{n}$$

viimasest võrratusest saame, et

$$S_w^2 > S_y^2$$

ehk süstemaatiline valik on lihtsast juhuvalikust efektiivsem, kui tunnuse y valimisisene hajuvus on suur, võrreldes hajuvusega ÜK-s.

Parim hinnang saadakse loendi korral, mis on järjestatud uuritava tunnuse või sellega tugevalt korreleeritud tunnuse väärtuste järgi.

Halva järjestusega loendi puhul võidakse valimisse saada liialt vähe varieeruvad objektid, mille tagajärjeks on ebatäpsed hinnangud alahinnatud usaldusintervalliga.

6.3 SÜ realiseerimine praktikas

1. Mõnikord on SÜ probleemiks, et pole võimalik saavutada täpselt etteantud valimimahtu. Näiteks kui $N = 125$ ja samm $m = 3$, saame $n = \left[\frac{125}{3}\right] = 41$ ehk valimimaht on kas 41 või 42 sõltuvalt juhuslikust stardist. Kui aga $m = 4$, siis $n = 31$, $n + 1 = 32$. Valimimahtusid nt. 33 – 40 pole võimalik saada. Suurte üldkogumite korral see probleem kaob.

2. Valimimahu reguleerimiseks, kasutatakse teisi SÜ protseduure, millest üks on näiteks ringsüstemaatiline valik. Selle meetodi korral vaadeldakse loendit ringina, kus viimasele elemendile järgneb jälle esimene. Genereeritakse juhuslik arv $1 \leq r \leq N$ ja võetakse talle vastav objekt ning sammu m tagant iga järgnev objekt, kuni soovitud valimimaht on käes.

7 Ebavõrdsete tõenäosustega valik

Andes üldkogumi elementidele erinevaid kaasamistõenäosusi, on võimalik parandada leitavate hinnangute omadusi. Vaatame siinkohal lähemalt üht enimkasutatavamat viisi.

Olgu $t = \sum_U y_i$ ja sellele vastav nihketa hinnang $\hat{t} = \sum_U \frac{I_i y_i}{EI_i}$.

Kui valida disain nii, et oodatavad valikute arvud on võrdelised y väärtustele, ehk

$$EI_i \propto y_i \quad (EI_i = cy_i),$$

siis $y_i/EI_i \equiv c$ ja hinnang saab järgmist kuju:

$$\hat{t} = c \sum_U I_i.$$

Kui lisaks disain on fikseeritud mahuga, siis $\sum_U I_i = n$ ja hinnang lihtsustub veelgi rohkem:

$$\hat{t} = cn.$$

Võtame viimases avaldises mõlemalt pool keskväärtuse ja saame, et $c = t/n$. Järelikult, fikseeritud mahuga disainide korral iga kogusumma hinnang \hat{t} annab meile tulemuseks täpse parameetri t sõltumata realiseerunust valimist.

Juhusliku valimimahuga disainide korral $\sum_U I_i = \mathbf{n}$. Olgu n_s realiseerunud valimimaht. Sel juhul:

$$\begin{aligned} t &= E(\hat{t}) = c E\left(\sum_U I_i\right) = c E(\mathbf{n}) \\ \Rightarrow c &= \frac{t}{E(\mathbf{n})}, \end{aligned}$$

ja kogusumma hinnang:

$$\hat{t} = \frac{t}{E(\mathbf{n})} n_s.$$

Ülalpool kirjeldatud valikut nimetatakse suurusega võrdelise tõenäosusega valikuks (*Sampling with Probabilities Proportional to Size*, PPS).

Probleemid seotud PPS realiseerimisega:

- kuna y_i pole teada enne valimi võtmist, siis ka pole võimalik leida $EI_i \propto y_i$;
- juhul, kui on võimalik kasutada taustinfot, ütleme tunnust x , mis on teadaolevalt positiivselt seotud uuritava tunnusega, siis saab valida $EI_i \propto x_i$;
- suurtes uuringutes, kus uuritavaid tunnuseid on palju, võib juhtuda, et EI_i on võrdelised ainult mõnede tunnustega; sellisel juhul hinnangud teiste tunnuste jaoks tulevad ebatäpsed.

PPS kasutamise **näiteid**...

1. Leibkonna eelarve uuring. Selleks kasutatakse tavaliselt rahvastikuregister (mis sisaldab infot inimeste kohta). Sellest võetakse valim võrdse tõenäosusega iga inimese jaoks. Leibkondadel on sellisel juhul tõenäosus olla valitud võrdeline leibkonna suurusega. Selline valikuviiis suurendab hinnangute täpsust, mis on seotud näiteks kulutustega, kuna need tunnused on enamasti tugevalt ja ka positiivselt korreleeritud leibkonna suurusega. Tulud on samuti positiivselt seotud leibkonna suurusega, kuid see seos on nõrgem.

2. Kui tahetakse hinnata vabade töökohtade arvu linnas, siis LJV puhul on valimis enamik väikeettevõtteid (neid on rohkem), aga hinnatav parameeter oleneb just palju suurfirmadest. Seega peaks neil olema suurem võimalus valimisse sattuda.

PPS kasutamise pealmised **põhjused**:

1. hinnangute täpsuse suurendamine;
2. kindlate objektide sattumine valimisse (nt nende edaspidiseks uurimiseks).

Märkus. Miks kutsutakse antud valikut 'ebavõrdsete tõenäosustega valikuks'?

kui tingimused on määratud keskväärtustele EI_i . Kus on siin tõenäosused?

Teame, et

$$EI_i = \begin{cases} \pi_i, & \text{TTA disainide korral;} \\ n p_i, & \text{TGA disainide korral.} \end{cases}$$

Järelikult, tingimused EI_i jaoks tähendavad ka tingimusi tõenäosuste jaoks.

7.1 Suurusega võrdelise tõenäosusega valik

Eeldame, et enne uuringu teostamist teame tausttunnuse x väärtuseid. Täpsemalt on selliseks tunnuseks mingit suurust iseloomustav tunnus.

Disaini moodustamiseks valime

$$EI_i \propto x_i,$$

mis tähendab, et

$$EI_i = c x_i \cdot \left| \sum_{i=1}^N (\dots) \right|$$

$$\underbrace{\sum_{i=1}^N EI_i}_{E\mathbf{n}} = c \underbrace{\sum_{i=1}^N x_i}_{t_x}$$

Järelikult,

$$c = \frac{E\mathbf{n}}{t_x}.$$

Kokkuvõttes võib öelda, et valikuindikaatori keskväärtus peab olema:

$$EI_i = \begin{cases} E(\mathbf{n})x_i/t_x, & \text{juhusliku valimimahuga disainide korral;} \\ n x_i/t_x, & \text{fikseeritud mahuga disainide korral.} \end{cases}$$

Kuigi valemid näevad lihtsad välja, pole siiski lihtne konstrueerida algoritmi fikseeritud mahuga TTA valiku teostamiseks. Üks tuntumaid on nn Sunteri algoritm (vaatame loengul, tahvlil).

TGA valik fikseeritud mahuga ei ole midagi muud kui multinomiaalne disain valikutõenäosustega $p_i = x_i/t_x$ ja valimimahuga n . Teostamisviisi vaatasime varem. Probleemiks - TGA disainid pole kõige eelistatumad disainid praktikas.

Üks lihtsamatest TTA disainidest on Poissoni valik, mis kahjuks annab juhuslikku valimimahtu. Kuid oma lihtsuse tõttu see disain on üsna populaarne praktikas.

7.2 Poissoni valik

Poissoni valiku korral kõik elemendid läbitakse järjest, alates esimesest kuni viimaseni, üks kord. Iga elemendi jaoks saadakse juhusliku valikuindikatori realisatsioon, $I_i \sim Be(\pi_i) = Bin(1, \pi_i)$, I_i on kõik sõltumatud juhuslikud suurused.

PPS valiku korral $\pi_i = nx_i/t_x$, x on tausttunnus. Meeldetuletuseks, Poissoni disaini karakteristikud:

$$\begin{aligned} EI_i &= \pi_i = nx_i/t_x, \\ VI_i &= \pi_i(1 - \pi_i), \\ E(I_i I_j) &= \pi_{ij} \overset{I_i \perp I_j}{=} \pi_i \pi_j, \\ Cov(I_i, I_j) &= 0. \end{aligned}$$

Teoreem 8. Poissoni valiku korral, $I \sim Bin(1, \pi_i)$, nihketa hinnang ÜK kogusummale $t = \sum_U y_i$ on järgmine:

$$\hat{t} = \sum_U \frac{I_i y_i}{\pi_i},$$

vastava punktihinnanguga:

$$\hat{t} = \sum_s \frac{y_i}{\pi_i}.$$

Hinnangu dispersioon on $V(\hat{t}) = \sum_U \frac{1-\pi_i}{\pi_i} y_i^2$

ja dispersiooni hinnang: $\hat{V}(\hat{t}) = \sum_U \frac{1-\pi_i}{\pi_i^2} y_i^2 I_i$
vastava punktihinnanguga

$$\hat{V}(\hat{t}) = \sum_s \frac{1-\pi_i}{\pi_i^2} y_i^2.$$

Tõestada teoreemi 8 väiteid iseseisvalt, arvestades, et $\Delta_{ij} = 0, i \neq j; \Delta_{ii} = VI_i, \pi_{ii} = \pi_i$.

Kuna Poissoni valik on juhusliku valimimahuga disain, siis eelistatakse alternatiivset (suhtetüüpi) hinnangut ÜK kogusummale:

$$\hat{t}_{alt} = \frac{\hat{t}}{\hat{N}} N,$$

kus $\hat{N} = \sum_s 1/\pi_i$.

8 Kihtvalik

Kihtvalik on praktikas enim kasutatav valikudisain, mille korral jaotatakse objektid ÜK-s mõne tausttunnuse (*kihistava tunnuse*) väärtuse järgi osadesse (*kihtidesse*). Kihte vaadeldakse üksteisest sõltumatute kogumitena, milledes võib rakendada erinevaid valikumeetodeid.

Kihtvalikut kasutatakse:

- Hinnangu täpsuse tõstmiseks - uuritava tunnuse suhtes homogeensed kihid ($y_i/EI_i \approx const$ ehk $y_i \propto EI_i$) tagavad valimihinnangu väikese varieeruvuse (disainiefekt < 1).
- Osakogumite hindamiseks - eriti väikeste valimite korral on mõttekas osakogumit esitada eraldi kihtidena, et rakendada seal temale sobivat optimaalset disaini.
- Erinevat käsitlust vajavate kihtide hindamine - kallimalt uuritavate objektide valimit vähendatakse, suure kao korral valimit suurendatakse.
- Uuringu administreerimine - suunamaks valimi paigutust (nt. intervjuerijate keskuste ümber). See võimaldab kokkuvõidu uuringu läbiviimisel.

Kihistava(te) tunnus(t)e valik:

- määratud üldkogumi kõigil objektidel, teada enne uuringu läbiviimist (sugu, vanus, maakond, linn/maa, töötajate arv,...)
- ei määra liiga peent kihistust, mis raskendaks osakogumite hinnangute leidmist.

Disain

Olgu lõplik üldkogum $U = \{1, \dots, N\}$ jagatud H kihiks $U_1, \dots, U_h, \dots, U_H$ vastavate mahtudega $N_1, \dots, N_h, \dots, N_H$ kihtides, kusjuures

$$U = \bigcup_{h=1}^H U_h, \quad U_h \cap U_g = \emptyset \text{ kui } h \neq g,$$

$$N = \sum_{h=1}^H N_h.$$

Tähistame valikuvektorit kihis h : $\mathbf{I}_h = (I_r, \dots, I_{r+N_h})$, kus r on eelmiste kihtide objektide arv + 1, $r = \sum_{i=1}^{h-1} N_i + 1$. Igas kihis rakendatakse teiste kihtide omast sõltumatut valikut vastavalt disainile $p_h(\mathbf{k}_h) = P(\mathbf{I}_h = \mathbf{k}_h)$.

Terve valikuvektor \mathbf{I} koosneb kihtide alamvektoritest,

$$\mathbf{I} = (\mathbf{I}_1, \dots, \mathbf{I}_h, \dots, \mathbf{I}_H),$$

ning tänu alamvektorite sõltumatusele saab valikudisaini esitada kihtide disainide korrutisena:

$$p(\mathbf{k}) = \prod_{h=1}^H p_h(\mathbf{k}_h),$$

kus $\mathbf{k} = (\mathbf{k}_1, \dots, \mathbf{k}_H)$.

8.1 Hindamine kihtvaliku korral

Tähistame:

$t_h = \sum_{U_h} y_i$ - uuritava tunnuse summa kihis U_h ,
 $\bar{Y}_h = \frac{t_h}{N_h}$ - keskmine kihis U_h .

Soovime hinnata ÜK kogusumma t ,

$$t = \sum_{h=1}^H t_h,$$

ehk alternatiivselt,

$$t = \sum_{h=1}^H N_h \bar{Y}_h.$$

Hinnatavaks parameetriks võib olla ka ÜK keskmine

$$\bar{Y} = \frac{t}{N} = \sum_{h=1}^H \frac{N_h}{N} \bar{Y}_h = \sum_{h=1}^H W_h \bar{Y}_h,$$

kus W_h on kihi osakaal ÜK-s.

Teoreem 9 (Kihtvalik). Kihtvaliku korral on nihketa hinnang ÜK summale t järgmine:

$$\hat{t} = \sum_{h=1}^H \hat{t}_h,$$

kus $E(\hat{t}_h) = t_h$ ehk hinnang \hat{t}_h on nihketa kihis U_h . Hinnangu \hat{t} dispersioon on

$$V(\hat{t}) = \sum_{h=1}^H V(\hat{t}_h)$$

ja sellele vastav nihketa hinnang

$$\hat{V}(\hat{t}) = \sum_{h=1}^H \hat{V}(\hat{t}_h),$$

kus $E[\hat{V}(\hat{t}_h)] = V(\hat{t}_h)$.

Teoreemi tõestus järeldeb hinnangute \hat{t}_h sõltumatuses erikihtides, samuti ka operaatorite E ja V omadustest.

Järeldus. Kihtvaliku korral avaldub hinnang ÜK keskmisele kihikeskmiste hinnangute kaalutud keskmisena,

$$\hat{Y} = \sum_{h=1}^H W_h \hat{Y}_h,$$

mille dispersioon on

$$V(\hat{Y}) = \sum_{h=1}^H W_h^2 V(\hat{Y}_h).$$

Kui kihtides kasutatakse nihketa hinnanguid dispersioonidele $\hat{V}(\hat{Y}_h)$, siis nihketa hinnang dispersioonile on

$$\hat{V}(\hat{Y}) = \sum_{h=1}^H W_h^2 \hat{V}(\hat{Y}_h).$$

8.2 Lihtne juhuslik kihtvalik

Kui kõikides kihtides kasutatakse lihtsat juhuvalikut TTA, siis nimetatakse sellist valikumeetodit *lihtsaks juhuvalikuks* (LJKV). Seejuures võib kihtides kasutada erinevaid valikusuhteid

$$f_h = \frac{n_h}{N_h}, h = 1, \dots, H.$$

Paneme tähele, et kuigi ühe kihi piires disain on isekaaluv, pole ta seda terves üldkogumis, mille tõttu valimikeskmine ja osakaal ei ole nihketa hinnangu- teks ÜK keskmisele ja osakaalule.

LJ TTA korral on kihi sees hinnanguks prameetrile t_h

$$\hat{t}_h = \sum_{U_h} \frac{I_i y_i}{\pi_i} = \frac{N_h}{n_h} \sum_{U_h} I_i y_i,$$

või valimi kaudu:

$$\hat{t}_h = N_h \bar{y}_h,$$

kus $\bar{y}_h = \frac{1}{n_h} \sum_{s_h} y_i$ valimikeskmine kihis U_h . Kasutades teoreemi (Kihtvalik ja LJ valik TTA) saame sõnastada teoreemi LJKV jaoks.

Teoreem (Lihtne juhuslik kihtvalik). Lihtsa juhusliku kihtvaliku korral avaldub kogusumma $t = \sum_U y_i$ kujul

$$\hat{t} = \sum_{h=1}^H N_h \bar{y}_h,$$

dispersiooniga

$$V(\hat{t}) = \sum_{h=1}^H N_h^2 (1 - f_h) S_{yh}^2 / n_h$$

ja dispersiooni nihketa hinnanguga

$$\hat{V}(\hat{t}) = \sum_{h=1}^H N_h^2 (1 - f_h) s_{yh}^2 / n_h,$$

kus

$$S_{yh}^2 = \frac{1}{N_h - 1} \sum_{U_h} (y_i - \bar{Y}_h)^2,$$

$$s_{yh}^2 = \frac{1}{n_h - 1} \sum_{s_h} (y_i - \bar{y}_h)^2.$$

Järeldus. Arvestades seost $\bar{Y} = \frac{t}{N}$, avalduvad vastavad avaldised **keskmise hinnangu** puhul järgmiselt:

$$\hat{\bar{Y}} = \frac{1}{N} \sum_{h=1}^H N_h \bar{y}_h,$$

$$V(\hat{\bar{Y}}) = \frac{1}{N^2} \sum_{h=1}^H N_h^2 (1 - f_h) S_{yh}^2 / n_h,$$

$$\hat{V}(\hat{\bar{Y}}) = \frac{1}{N^2} \sum_{h=1}^H N_h^2 (1 - f_h) s_{yh}^2 / n_h.$$

Näide 1. Reklaamifirmat huvitab, kui palju teha reklaami ühes maakonnas, kuhu kuulub 2 suuremat linna (A ja B) ning maapiirkond. Selleks uuritakse, mitu tundi nädalas inimesed maakonnas keskmiselt televiisorit vaatavad.

- U_1 : Linn A on ehitatud suure tehase juurde ning enamiku linna elanikkonnast (155 majapidamist) moodustavad tehase töötajad kooliealiste lastega.
- U_2 : Linn B on naabruses asuva suure linna eeslinnaks ning enamik 62 majapidamisest on vanemad inimesed väheste lastega.
- U_3 : Maapiirkonnas elab 93 majapidamist.

Raha on 40 majapidamise küsitlemiseks. Otsustatakse, et valimid kihiti on $n_1 = 20$, $n_2 = 8$, $n_3 = 12$. Igast kihist võetakse LJ valim. Tulemused – TV ees veedetud tunnid nädalas – on toodud allolevas tabelis.

Kiht 1, Linn A	Kiht 2, Linn B	Kiht 3, Maa
35 28 26 41 43 29 32 37 36 25 29 31 39 38 40 45 28 27 35 34	27 4 49 10 15 41 25 30	8 15 21 7 14 30 20 11 12 32 34 24
$n_1 = 20$ $\bar{y}_1 = 33,9$ $s_{y1}^2 = 35,358$ $N_1 = 155$	$n_2 = 8$ $\bar{y}_2 = 25,125$ $s_{y2}^2 = 232,411$ $N_2 = 62$	$n_3 = 12$ $\bar{y}_3 = 19$ $s_{y3}^2 = 87,636$ $N_3 = 93$

Leiame hinnangu TV vaatamisele nädakeskmisele majapidamise kohta koos usalduspiiridega ning suhtelise veaga.

$$\hat{Y} = \frac{1}{N} \sum_{h=1}^H N_h \bar{y}_h = \frac{1}{310} (155 \cdot 33,9 + 62 \cdot 25,125 + 93 \cdot 19) = 27,7;$$

$$\hat{V}(\hat{Y}) = \frac{1}{N^2} \sum_{h=1}^H N_h^2 (1-f_h) s_{yh}^2 / n_h = \frac{1}{310^2} \left[\frac{155^2 \cdot 0,871 \cdot 35,358}{20} + \frac{62^2 \cdot 0,871 \cdot 232,411}{8} + \frac{93^2 \cdot 0,871 \cdot 87,636}{12} \right] = 1,97.$$

Tõenäosusega 95% saame väita, et keskmiselt vaadatakse TV majapidamises $27,7 \pm 1,96 \cdot \sqrt{1,97} = 27,7 \pm 2,8$ tundi nädalas.

Punkti hinnangu suhteline viga, $Suht.v.(\hat{Y}) = \frac{\sqrt{\hat{V}(\hat{Y})}}{\hat{Y}} = \frac{\sqrt{1,97}}{27,7} = 5,1\%$, mis on väga kõrge näitaja ning järelikult, võib hinnangut usaldada.

Samas, paneme tähele, et kui leida hinnangud kihiti eraldi ning hinnangu täpsust, siis linnas B tuleb suhteline viga 20,04%, mis on äärmiselt suur.

Ülesanne. Leida linna B hinnang ning hinnangu suhteline viga koos hinnangu usalduspiiridega. Mis aitaks muuta hinnangut täpsemaks?

8.3 Valimi optimaalne paigutus

Kihtvaliku teostamisel on esmaseks ülesandeks kihtide moodustamine üldkogumis. Fikseeritakse tunnused, mille abil objektid jagatakse kihtidesse. Teiseks tähtsaks ülesandeks on valikudisainide määramine kihtides. Kolmandaks oluliseks ülesandeks on valimimahtude määramine kihtides.

Olgu kogu valimimaht n . Temast sõltub hinnangute täpsus. Mida suurem n , seda väiksem dispersioon. Samas mahtu n suurendades kasvab ka uuringu maksumus. Uuringu maksumus on tavaliselt eelarvega fikseeritud, mis paneb kitsendused ka valimimahule. Osutub aga, et valimimahtu n oskuslikult kihtidesse jagades võime nii hinnangute dispersioone kui ka uuringu maksu- must vähendada.

Olgu hinnangu \hat{t}_y dispersioon avaldav kujul

$$V = V(\hat{t}_y) = \sum_{h=1}^H \frac{A_h}{n_h} + B, \quad (34)$$

kus kihtidesisese hajuvus komponendid A_h ja üldine komponent B ei sõltu valimimahtudest n_h . Olgu uuringu kogumaksumus C avaldatav seosega

$$C = c_0 + \sum_{h=1}^H n_h c_h, \quad (35)$$

kus c_0 on üldkulud ja c_h on andmete saamise kulu h -nda kihi objektilt. Suurused c_0 ja c_h on uuringut planeerides teada.

Eesmärgiks on saavutada valimimahtude n_h määramisega kihtides parimaid tulemusi dispersiooni ja maksumuse seisukohalt.

Defnitsioon. Valimimahtude komplekti $n_h, h = 1, \dots, H$, nimetatakse optimaalseks, kui kehtib üks järgmistest tingimustest:

1. Etteantud uuringu kogumaksumuse C juures on hinnangu dispersioon $V = V(\hat{t}_y)$ minimaalne.
2. Etteantud hinnangu dispersiooni V juures on uuringu kogumaksumus minimaalne.
3. Etteantud valimimahu n juures on nii dispersioon kui ka maksumus minimaalsed.

Järgnevas tõestame teoreemi, mis annab optimaalsed valimimahud n_h kõigi ülalloeletud eesmärkide saavutamiseks.

Teoreem 11 (valimi optimaalsest paigutusest). Kihtvaliku korral, kus hinnangu dispersioon V ja maksumus C on antud valemitega (34)-(35), saavutatakse valimi optimaalne paigutus, kui

$$n_h \propto \sqrt{\frac{A_h}{c_h}}, h = 1, \dots, H. \quad (36)$$

Tõestus. Ülalloetletud optimaalsuse eesmärkide saavutamiseks tuleb minimeerida korrutis $V \cdot C$ suuruste n_h suhtes. Jättes kõrvale suurustest n_h mittesõltuvad liikmed, tuleb minimeerida korrutis

$$K = \left(\sum_{h=1}^H \frac{A_h}{n_h} \right) \left(\sum_{h=1}^H n_h c_h \right).$$

Kasutame Cauchy-Schwarzi võrratust

$$\sum a_i^2 \sum b_i^2 \geq \left(\sum a_i b_i \right)^2,$$

kus võrdus kehtib parajasti siis, kui $b_i/a_i = \text{const}, \forall i$. Saame, et mistahes n_h valiku korral

$$K \geq \left(\sum_{h=1}^H \sqrt{A_h c_h} \right)^2.$$

Kuna parem pool ei sõltu suurustest n_h , siis selline n_h valik, mis annab võrduse annab ka K minimaalse väärtuse. Cauchy-Schwarzi teoreemist saame, et võrdus kehtib kui

$$n_h \sqrt{\frac{c_h}{A_h}} = \text{const} \text{ ehk } n_h \propto \sqrt{\frac{A_h}{c_h}}.$$

Sellega on teoreem tõestatud. \diamond

Näeme, et valimi optimaalseks paigutamiseks kihtidesse tuleb rohkem objekte valida sellest kihist, kus kihisisene dispersioonikomponent A_h on suur, aga maksumus c_h väike. Võrdeteguri leidmine sõltub püstitatud optiseerimisülesandest.

Teoreem 12. Dispersiooni V minimiseerib fikseeritud maksumuse C korral järgmine valimi paigutus

$$n_h = (C - c_0) \frac{\sqrt{A_h/c_h}}{\sum_{h=1}^H \sqrt{A_h c_h}}, h = 1, \dots, H, \quad (37)$$

ja minimaalne dispersioon on

$$V_{opt} = \frac{1}{C - c_0} \left(\sum_{h=1}^H \sqrt{A_h c_h} \right)^2 + B. \quad (38)$$

Tõestus. Seosest (36) järeldub, et mingi konstandi λ korral kehtib

$$n_h = \lambda \sqrt{\frac{A_h}{c_h}}, h = 1, \dots, H.$$

Asendades saadud seose maksumuse avaldise (35), saame võrdeteguri λ jaoks,

$$\lambda = \frac{C - c_0}{\sum_{h=1}^H \sqrt{A_h c_h}}.$$

Viimased kaks seost annavadki teoreemi väite (37) n_h kohta. Kasutades optimaalseid valimimahte dispersiooniavaldises (34), saame teoreemi väite (38).

\diamond

Teoreem 13. Maksimuse C minimiseerib fikseeritud dispersiooni V korral järgmine n_h planeering,

$$n_h = \sqrt{\frac{A_h}{c_h}} \cdot \frac{\sum_{h=1}^H \sqrt{A_h c_h}}{V - B}, h = 1, \dots, H, \quad (39)$$

ja vastav optimaalne maksimum sel juhul on

$$C_{opt} = c_0 + \frac{1}{V - B} \left(\sum_{h=1}^H \sqrt{A_h c_h} \right)^2. \quad (40)$$

Tõestus. Analoozne eelmisele teoreemile. \diamond

Kogu valimimaht optimaalsete kihisestest valimimahtude korral on $n = \sum_{h=1}^H n_h$.

8.4 Optimaalne valimi paigutus KLJV korral

Kiit lihtne juhuslik valik on praktikas sageli kasutatav disain. Teame, et kogusumma hinnangu dispersioon avaldub sel juhul,

$$V(\hat{t}_y) = \sum_{h=1}^H \frac{N_h^2}{n_h} (1 - f_h) S_{yU_h}^2.$$

Valemist näeme, et see dispersioon avaldub just nii nagu meie tulemusteks vaja:

$$A_h = N_h^2 S_{yU_h}^2, B = - \sum_{h=1}^H N_h S_{yU_h}^2.$$

Teoreem valimi optimaalsest paigutusest ütleb nüüd, et

$$n_h \propto \frac{N_h S_{yU_h}}{\sqrt{c_h}}. \quad (41)$$

Näeme, et valimi optimaalseks planeerimiseks peame võtma rohkem objekte kihist, mille maht N_h on suurem, milles tunnuse y dispersioon on suurem, aga milles objekti küsitlemine/mõõtmine on odavam. Fikseeritud maksimumuse korral on optimaalseks planeeringuks,

$$n_h = \frac{C - c_0}{\sum_{h=1}^H N_h S_{yU_h} \sqrt{c_h}} \cdot \frac{N_h S_{yU_h}}{\sqrt{c_h}}. \quad (42)$$

8.5 Alternatiivsed valimi paigutused KLJV korral

Olgu nüüd $c_h = c(const), \forall h$. Tänapäeva praktikas on see enamasti toimiv eeldus. Maksumuse avaldisest (35) saame nüüd, et

$$C - c_0 = c \cdot n. \quad (43)$$

Seega kui uuringu kogumaksumus on ette antud, on sellega fikseeritud ka kogu valimimaht n .

1. Neymani paigutus (1934). Valimimahtude Neymani paigutus on optimaalne paigutus (42) fikseeritud maksumuse korral, kui $c_h = const$. Siis saame valemist (42)-(43) erijuhu,

$$n_h = n \frac{N_h S_{yU_h}}{\sum_{h=1}^H N_h S_{yU_h}}. \quad (44)$$

Paneme tähele, et Neymani paigutus, nagu ka kõik eelnevad valimi paigutused on optimaalsed tunnuse y jaoks. Mõne teise tunnuse z jaoks ei pruugi selline valimi jagamine hea olla.

Näide. Koosnegu üldkogum kolmest kihist mahtudega $N_1 = 150, N_2 = 90, N_3 = 120$. Eelmistest uuringutest on teada, et $S_{yU_1} = 100, S_{yU_2} = 200, S_{yU_3} = 300$. Eeldades konstantset maksumust saaksime optimaalseks paigutuseks kogu valimimahu 12 korral $n_1 = 2.6, n_2 = 3.1, n_3 = 6.3$, ehk ümardatult $n_1 = 3, n_2 = 3, n_3 = 6$.

2. Võrdeline paigutus. Sel juhul on vastavate kihtide osakaalud valimis ja üldkogumis võrdsed:

$$n_h = n \frac{N_h}{N}. \quad (45)$$

Sel juhul on valikusuhted kihtides võrdsed: $f_h = n_h/N_h = n/N = f$. Valemist (44) näeme, et selline paigutus on optimaalne, kui uuritava tunnuse dispersioonid on kõigis kihtides võrdsed, muidugi ka $c_h = const, \forall h$.

Näeme, et võrdeline planeering on tunnuse iseloomu suhtes neutraalne, ühtviisi hea kõikide tunnuste jaoks, aga ei pruugi olla optimaalne ühegi tunnuse jaoks.

Näide. Võrdeline paigutus annab eelmise näite andmetel $n_1 = 5, n_2 = 3, n_3 = 4$.

3. x-optimaalne paigutus. Kuna uuritav tunnus ei ole enne uuringut teada, siis tehakse valimi paigutus kasutades temaga tugevasti korreleeritud teadaolevat x -tunnust.

4. Kogusummaga võrdeline paigutus. Olgu $t_y = \sum_U y_i$ ja $t_{yU_h} = \sum_{U_h} y_i$. Olgu $y_i \geq 0, \forall i$, siis

$$n_h = n \frac{t_{yU_h}}{t_y}.$$

See paigutus on optimaalne, kui variatsioonikordajad on kihiti võrdsed (veendu!):

$$CV_h = \frac{S_{yU_h}}{\bar{Y}_{U_h}} = \text{const}, \forall h.$$

8.6 LJV ja KLJV võrdlemine

Kogusumma nihketa hinnanguks on TTA disainide korral $\sum_s y_i / \pi_i$. Tahame võrrelda selle hinnangu dispersiooni LJV ja KLJV korral. LJV korral $f_i = n/N$ ja hinnang teiseneb kujule $\hat{t}_y = N\bar{y}$, tema dispersiooniks on

$$V_{LJV}(\hat{t}_y) = \frac{N^2}{n}(1-f)S_y^2. \quad (46)$$

KLJV korral on $\pi_i = n_h/N_h$ kui $i \in U_h$ ja nihketa hinnang saab kuju $\hat{t}_y = \sum_{h=1}^H N_h \bar{y}_h$. Selle dispersioon on

$$V_{KLJV}(\hat{t}_y) = \sum_{h=1}^H \frac{N_h^2}{n_h}(1-f_h)S_{yU_h}^2. \quad (47)$$

Kumma disaini korral on hinnang on täpsem, kui kogu valimimaht n on sama? Sõltub paljudest asjaoludest. Valimi õige planeerimisega on võimalik saavutada antud kihistuse korral minimaalne dispersioon. Kui maksumus on sama, siis dispersiooni minimiseerib Neymani planeering ja vastav optimaalne dispersioon on

$$V_{opt}(\hat{t}_y) = \frac{N^2}{n} A \sum_{h=1}^H W_h (1-f_h) S_{yU_h}^2, \quad (48)$$

kus $W_h = N_h/N$ on kihi osakaal ja $A = \sum_{h=1}^H W_h S_{yU_h}$. Kui dispersioonid kihtides on võrdsed $S_{yU_h} = S_{y0}$ teiseneb dispersioonivalem eriti lihtsale kujule:

$$V_{opt}(\hat{t}_y) = \frac{N^2}{n}(1-f)S_{y0}^2. \quad (49)$$

Kui me ei tea midagi kihi dispersioonidest arvata, kuid kasutame valimi võrdelist paigutust, saame hinnangu dispersioonile valemist (47) kuju:

$$V_{vord}(\hat{t}_y) = N^2(1-f) \sum_{h=1}^H W_h^2 S_{yU_h}^2 / n_h. \quad (50)$$

Samas, kui dispersioonid kihtides juhtuvad olema võrdsed, annab see valem sama dispersiooni, mis optimaalne valem (49).

Ül. Näita see väide. Samuti tuleta valemid (49)-(50).

Valemist (49) näeme taaskord üht kihtvaliku printsiipi, kui objektid on kihtidesse jagatud selliselt, et S_{y0}^2 on väike, on ka hinnangu \hat{t}_y dispersioon väike. Põhimõtteliselt võib kihistamisega saavutada nulldispersiooni.

Üldjuhul, kui valimimaht on sama, kehtivad dispersionide vahel järgmised seosed:

$$V_{opt}(\hat{t}_y) \leq V_{vord}(\hat{t}_y) \leq V_{LJV}(\hat{t}_y).$$

Kokkuvõtteks. KLJV kasutamine LJV asemel on hinnangute täpsuse seisukohalt õigustatud, kui

1. Tunnused on kihtide sees homogeensed (sarnased objektid on samas kihis).
2. Tunnuste keskmised on kihiti erinevad.

Valimimahu võrdeline paigutus on hea, optimaalne paigutus annab väga hea tulemuse kindla uuritava tunnuse korral. Suuremahulistes uuringutes, kus uuritavaid tunnuseid on palju, on mõttekas kasutada võrdelist paigutust, et saada võimalikult hea hinnang kõigi tunnuste korral.

Näide. Olgu üldkogum mahuga $N = 6$ jagatud kaheks kihiks, nii et esimesed 3 objekti ühes ja järgmised kolm teises kihis. Seega $N_1 = 3$ ja $N_2 = 3$.

Olgu teada ka uuritava tunnuse väärtused $y = (2, 0, 1, 5, 9, 4)$. Olgu $n = 4$. Võrdleme LJV ja võrdelise planeeringuga KLJV, st $n_1 = n_2 = 2$.

Näeme, et üldkogumis $\bar{Y} = 3, 5$; $\bar{Y}_1 = 1$; $\bar{Y}_2 = 6$ ja $S_y^2 = 10, 7$; $S_{yU_1}^2 = 1$, $S_{yU_2}^2 = 7$. Keskvärtuse hinnanguks on LJV korral valimikeskmine $\bar{y} = \sum_s y_i/4$. Võrdelise planeeringuga KLJV korral tuleb selleks samuti tavaline valimikeskmine. Leiame valimikeskmise dispersioonid:

$$V_{LJV}(\hat{Y}) = (1 - f)S_y^2/n = \frac{6 - 4}{6} \cdot \frac{10,7}{4} = 0,89$$

$$V_{KLJV}(\hat{Y}) = \sum_{h=1}^2 W_h^2(1-f)S_{yU_h}^2/n_h = \left(\frac{3}{6}\right)^2 \frac{3-2}{3} \cdot \frac{1}{2} + \left(\frac{3}{6}\right)^2 \frac{3-2}{3} \cdot \frac{7}{2} = 0,33.$$

Kommenteeri, mis aitas kaasa dispersiooni vähenemisele.

9 Järelkihistamine

Järelkihistamine on hinnangute täpsuse tõstmise meetod. Seda teostatakse hinnangute arvutamise etapil, st siis kui andmed valimilt on juba kogutud. Seejures valim võib olla võetud mistahes valikudisainiga.

Järelkihistamisel jagatakse valimi objektid gruppidesse – järelkihtidesse. Selleks peab valimi objektidel olema mõõdetud tunnus(ed), mida järelkihistamisel kasutatakse. Üldkogumi tasemel on vaja teada järelkihtide mahtusid. Kui vajalikud suurused on teada, saab moodustada mitmesuguseid järelkihistusi.

Järelkihistamist kasutatakse hinnangute täpsuse tõstmiseks. Kui õnnestub valim jagada gruppidesse nii, et objektid nendes on võimalikult homogeen- sed, siis väheneb hinnangute dispersioon.

Järelkihistamist saab kasutada ka kaost põhjustatud nihke vähendamiseks. Selle saavutamiseks jagatakse vastanute valim järelkihtidesse nii, et nendes vastanud on sarnased mittevastanutega.

Järelkihid on oma olemuselt osakogumid. Valimimaht nendes on juhuslik. Üldjuhul pole aga eesmärgiks hinnangute leidmine nendes osakogumites, neid

kasutatakse üldkogumihinnangute täpsustamiseks.

Olgu H järelkihti U_h . Need on mittelõikuvad ja ammendavad üldkogumi osad,

$$U = \bigcup_{h=1}^H U_h.$$

Olgu N_h järelkihi maht ja \bar{Y}_h järelkihi keskmine:

$$\bar{Y}_h = \frac{\sum_{U_h} y_i}{N_h}.$$

Olgu üldkogumist võetu valim s , mille osa järelkihis U_h tähistame s_h . Olgu disainikaalud $w_i = I_i/EI_i$. Järelkihi keskmise hinnanguks võtame suhte tüüpi hinnangu

$$\tilde{y}_h = \frac{\hat{t}_{yh}}{\hat{N}_h}, \quad (51)$$

kus $\hat{t}_{yh} = \sum_{s_h} w_i y_i$ ja $\hat{N}_h = \sum_{s_h} w_i$. Keskmise (51) baasil moodustatud kogusumma hinnanguks järelkihis on

$$\hat{t}_{yh} = N_h \tilde{y}_h. \quad (52)$$

Näeme, et siin läheb vaja teada järelkihtide mahtusid. Järelkihthinnanguks üldkogumi kogusummale on

$$\hat{t}_{jarel} = \sum_{h=1}^H N_h \tilde{y}_h. \quad (53)$$

Järelkihthinnangu dispersiooni leidmine on keeruline, sest liidetavad y_h ei ole sõltumatud, nagu nad olid seda kihtvaliku korral. Hinnang \hat{t}_{jarel} on aga vaadeldav üldisemate regressioon ja kalibreerimishinnangute erijuhuna. Nende dispersiooniavaldised on teada (Särndal jt 1992). Siin toome dispersioonihinnangu valemi,

$$\hat{V}(\hat{t}_{jarel}) = \sum_s \sum \check{\Delta}_{ij} w_i e_i w_j e_j, \quad \text{kus } e_i = \frac{N_h}{\hat{N}_h} (y_i - \tilde{y}_i), i \in s_h. \quad (54)$$

9.1 Järelikihthinnang LJV korral

LJV korral on $w_i = N/n$, millest

$$\hat{N}_h = \sum_{s_h} w_i = \sum_{s_h} \frac{N}{n} = N \frac{n_h}{n}.$$

Analoogiliselt saame, et

$$\hat{t}_{yh} = \sum_{s_h} w_i y_i = \frac{N}{n} \sum_{s_h} y_i.$$

Avaldisest (51) ja viimasest kahest võrrandist saame kokku LJV korral

$$\tilde{y}_h = \frac{1}{n_h} \sum_{s_h} y_i = \bar{y}_h.$$

Järelikihthinnang LJV korral avaldub valemist (53) järgmiselt:

$$\hat{t}_{jarel} = \sum_{h=1}^H N_h \bar{y}_h. \quad (55)$$

Saab näidata, et järelikihistus LJV korral on sama täpne eelkihistusega (KLJV) võrdelise planeeringu korral s.t.

$$V(\hat{t}_{jarel}) = N^2 \frac{1-f}{n} \sum_{h=1}^H \frac{N_h - 1}{N - 1} S_{yU_h}^2$$

ja

$$\hat{V}(\hat{t}_{jarel}) = N^2 \frac{1-f}{n} \sum_{h=1}^H \frac{N_h^2}{n_h} S_{y^{s_h}}^2.$$

10 Klastervalik

Vt. tund *Klastervalik* veebikeskkonnas Moodle.

11 Kahe-astmeline valik

Kahe-astmeline valik on protseduur, kus

1. astmel moodustatakse klaster-valim vastavalt mõnele tõenäosuslikele disainile ja
2. astmel valitud klastritest võetakse omakorda valimid.

Siinjuures 1. ja 2. astme valikudisainid ei sõltu üksteisest (võivad olla samad, võivad olla aga erinevad). Samuti ka eriklastrites võib rakendada erinevat valikudisaini.

Seega, **kihtvalik** on kahe-astmelise disaini erijuht, kui 1. astmel toimub kõikne klastervalik.

Klastervalik on kaheastmelise valiku erijuht, kui 2. astmel toimub kõikne valik igas valitud klastris.

Kirjanduses nimetatakse sageli 1. astme valikuühikuid (ehk klastreid) PSU=*Primary Sampling Unit*; 2. astme ühikuid - SSU=*Secondary sampling units*.

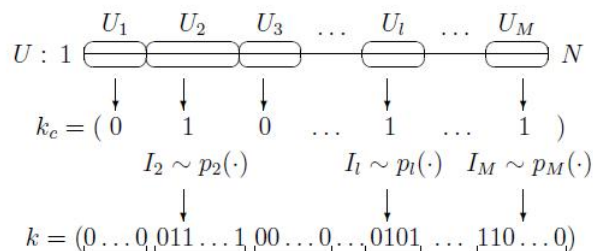
11.1 Tähistused

Olgu U_e - PSU, kusjuures $U_e \subset U$, $\bigcup_{e=1}^M U_e = U$, $U_e \cap U_g = \emptyset$ ja N_e klasteri U_e suurus.

Kasutame erinevaid valikuvektoreid:

$I_c = (I_{c1}, \dots, I_{ce}, \dots, I_{cM})$ valikuvektor klastrite (PSU-de) jaoks;
 $I_c \sim p_c(k_c)$ valikudisain PSU jaoks;
 $I_e \sim p_e(k_e), e = 1, \dots, M$ valikudisain SSU jaoks klasteri U_e sees;
 $I_e = (I_{\cdot|e}, \dots, I_{i|e}, \dots, I_{\cdot|e})$ valikuvektor pikkusega N_e SSU jaoks klastris U_e .
 Tähtsad eeldused:

- valikud erinevates klastrites (PSU-des) on üksteisest sõltumatud;
- valikud 2. astmel ei sõltu valikust 1. astmel.



Valimimahud:

- $m = \sum_{e=1}^M k_{ce}$ valitud klastrite arv = klastervalimi maht;
- $n = \sum_{i=1}^N k_i$ lõplik valimimaht.

11.2 Hindamine kahe-astmelise valiku korral

Olgu $Y_e = \sum_{i \in U_e} y_i$ uuritava tunnuse kogusumma klastris U_e . Siis ÜK kogusummat saab esitada kujul

$$t = \sum_{e=1}^M Y_e. \quad (56)$$

Kui klastrid on valitud, siis saame nendes leida nihketa hinnangud \hat{Y}_e klastri summadele. Kogu ÜK summat saab hinnata järgmise nihketa hinnangu abil:

$$\hat{t} = \sum_{e=1}^M \frac{\hat{Y}_e}{E(I_{ce})} I_{ce}. \quad (57)$$

Edaspidi eeldame, et tegemist on TTA klastervalikuga, ehk $I_{ce} \in \{0, 1\}$. (TGA disaini vaatame hiljem eraldi.)

Klastersumma Y_e hinnangu saame tavalise nihketa hinnangu abil (kasutades Üldist Hindamisteoreemi):

$$\hat{Y}_e = \sum_{i \in U_e} \frac{y_i I_{i|e}}{E(I_{i|e})}, \quad (58)$$

kus $I_{i|e}$ on i -nda objekti valikuindikaator klastris U_e .

Üldisest hindamisteoreemist teame, et hinnangu \hat{Y}_e dispersioon avaldub järgmiselt:

$$V(\hat{Y}_e) = \sum_{i,j \in U_e} \Delta_{ij|e} \frac{y_i}{E(I_{i|e})} \frac{y_j}{E(I_{j|e})} \quad (59)$$

ja tema nihketa hinnang järgmise valemi abil:

$$\hat{V}(\hat{Y}_e) = \sum_{i,j \in U_e} \frac{\Delta_{ij|e}}{E(I_{i|e}I_{j|e})} \frac{y_i}{E(I_{i|e})} \frac{y_j}{E(I_{j|e})} I_{i|e}I_{j|e}. \quad (60)$$

Need dispersioonide valemid kehtivad konkreetse klatri sees. Kogu kaheastmelise protsessi jooksul tekkinud varieeruvust pole nii lihtne leida. Peame arvestama varieeruvuse nii 1. kui ka 2. astmel. Siin saame kasutada tingliku keskvaärtuse ja tingliku dispersiooni valemid:

$$E(\hat{t}) = E_{p_c} E(\hat{t}|I_c), \quad (61)$$

$$V(\hat{t}) = E_{p_c} V(\hat{t}|I_c) + V_{p_c} E(\hat{t}|I_c). \quad (62)$$

Valem (61) tähendab, et esmalt leiame keskvaärtuse igas klatri eraldi (II astme valikudisaini suhtes) ja seejärel keskmistame need klatri keskmised omakorda (leiame keskvaärtuse I astme disaini suhtes).

Arvestades valemid (57)-(58), kontrollime hinnangu nihketuse omaduse valemi (61) abil:

$$E(\hat{t}) = E_{p_c} E(\hat{t}|I_c) = E_{p_c} E\left(\sum_{e=1}^M \frac{\hat{Y}_e}{E(I_{ce})} I_{ce} | I_c\right) = E_{p_c} \left(\sum_{e=1}^M \frac{Y_e}{E(I_{ce})} I_{ce}\right) = \sum_{e=1}^M Y_e = t, \quad (63)$$

kus $E(\hat{Y}_e|I_c) = Y_e$.

Hinnangu \hat{t} kogu dispersiooni saamiseks leiame kõigepealt valemi (62) 2. liidetava. Selleks paneme tähele, et valemist (63)

$$E(\hat{t}|I_c) = \sum_{e=1}^M \frac{Y_e}{E(I_{ce})} I_{ce}, \quad (64)$$

mis on nihketa hinnanguks parameetrile t klasterdisaini $p_c(\cdot)$ suhtes. Sellele hinnangule saame rakendada ÜHT, et leida tema dispersiooni $V_{p_c}[E(\hat{t}|I_c)]$ (ja

seda disaini $p_c(\cdot)$ suhtes):

$$V_{p_c} E(\hat{t}|I_c) = V_{p_c} [E(\hat{t}|I_c)] = \sum_{e=1}^M \sum_{g=1}^M \Delta_{ceg} \frac{Y_e}{E(I_{ce})} \frac{Y_g}{E(I_{cg})} = V_1, \quad (65)$$

kus $\Delta_{ceg} = Cov(I_{ce}, I_{cg})$.

Nüüd, leiame valemi (62) 1. liidetava,

$$V(\hat{t}|I_c) = V \left(\sum_{e=1}^M \frac{\hat{Y}_e}{E(I_{ce})} I_{ce} | I_c \right) =$$

Arvestades, et 2. astmel toimub valik klastritest üksteisest sõltumata, siis saab dispersioonimärgiga V summa sisse minna,

$$= \sum_{e=1}^M V(\hat{Y}_e) \frac{I_{ce}}{[E(I_{ce})]^2},$$

kus $V(\hat{Y}_e|I_c) = V(\hat{Y}_e)$ on 2. astme valiku dispersioon, mis ei sõltu 1. astme valikust. Lisaks, $I_{ce}^2 = I_{ce}$ TTA disaini jaoks.

Järgmisena,

$$E_{p_c} V(\hat{t}|I_c) = E_{p_c} \left[\sum_{e=1}^M V(\hat{Y}_e) \frac{I_{ce}}{(E I_{ce})^2} \right] = \sum_{e=1}^M \frac{V(\hat{Y}_e)}{E(I_{ce})} = V_2. \quad (66)$$

Sõnastame eelnevat teoreemina.

Teoreem (Kahe-astmeline valik, TTA). Kahe-astmelise disaini korral nihketa hinnang ÜK summale on antud valemite (57)-(58) abil. Selle hinnangu dispersioon avaldub kui

$$V(\hat{t}) = V_1 + V_2,$$

kus V_1 on antud valemis (65) ja V_2 valemis (66).

11.3 Kahe-astmeline lihtne juhuslik valik

Selle valiku korral toimub 1. astmel LJ klastervalik TTA, kus

$$f_c = E(I_{ce}) = \frac{m}{M} \quad (67)$$

ja teisel astmel igast valitud klastrist võetakse omakorda LJ valik TTA, kus

$$f_e = \frac{n_e}{N_e}. \quad (68)$$

Nüüd, valemitest (57)-(58) saab leida nihketa hinnangu ÜK summale:

$$\hat{t} = \sum_{e=1}^M \frac{M}{m} I_{ce} \sum_{i \in U_e} \frac{N_e}{n_e} I_{ie} y_i. \quad (69)$$

Vastav punkt hinnang on

$$\hat{t} = \frac{M}{m} \sum_{e \in s_c} \frac{N_e}{n_e} \sum_{i \in s_e} y_i = \frac{M}{m} \sum_{s_c} N_e \bar{y}_e,$$

kus

s_c on klastervalim;

$i \in s_e$ summa üle valimi klastrist U_e ;

\bar{y}_e valimikeskmise valimis s_e .

Hinnangu ÜK summale saab kaalude abil kirja panna järgmiselt:

$$\hat{t} = \sum_{i \in s} w_i y_i,$$

kus $w_i = \frac{M N_e}{m n_e}$. Paneme tähele, et objektidel erinevatest klastritest on erinevad kaalud!

Hinnangu dispersiooni $V(\hat{t})$ saab välja kirjutada Teoreemist (Kahe-astmeline disain).

11.4 Isekaaluv kahe-astmeline valik

Praktikute lemmik on isekaaluv kahe-astmeline valik, kus lõplikelt valikuühikutel on võrdsed kaalud. Sel juhul valimi struktuur vastab ÜK struktuurile ja

valimi karakteristikud (keskmine, osakaal) on hinnanguteks vastavatele ÜK parameetritele.

Eeldame, et mõlemal astmel on teostatud TTA valik. Siis

$E(I_{ce}) = \pi_{ce}$ – e -nda klasteri kaasamistõenäosus; $E(I_{i|e}) = \pi_{i|e}$ – i -nda objekti kaasamistõenäosus klasteris U_e .

Valemitest (57)-(58) hinnang ÜK summale tuleb järgmine:

$$\hat{t} = \sum_{e=1}^M \sum_{i \in U_e} \frac{y_i}{\pi_{i|e} \pi_{ce}} I_{i|e} I_{ce}.$$

Disain on isekaaluv, kui

$$\pi_{i|e} \pi_{ce} = c(\text{const}), \forall i \in U.$$

Seda on võimalik saavutada kahel viisil:

1. Klasterite kaasamistõenäosused on võrdelised klasteri suurustega, $\pi_{ce} = N_e \frac{m}{N}$, $\forall e$, kusjuures kehtib

$$m = \sum_{e=1}^M \pi_{ce}.$$

2. astmel valitakse iga klasteri jaoks võrdne objektide arv valimisse nii, et $\pi_{i|e} = \frac{n_0}{N_e}$. Kokku saame, et

$$\pi_{ce} \pi_{i|e} = N_e \frac{m}{N} \frac{n_0}{N_e} = \frac{mn_0}{N} = \frac{n}{N}, \forall i \in U.$$

Siin n on lõplik valimimaht.

Sellise disaini korral kõikidel intervjueriatel on võrdne arv inimesi küsitlemiseks klasterites.

2. Esimesel astmel valitakse klasterid võrdse tõenäosusega,

$$f_c = \frac{m}{M}, \forall e.$$

Teisel astmel valitakse objektid võrdse tõenäosusega, st igas klastris on konstantne kaasamistõenäosus,

$$\pi_{i|e} = \frac{n_e}{N_e} = f, (const).$$

Sellisesl juhul

$$\pi_{ce}\pi_{i|e} = \frac{m}{M}f,$$

mis on võrne iga objekti jaoks.

Selliseks disainiks on LJ valik mõlemal astmel, kus 2. astmel toimub valimi võrdeline paigutus valitud klastritesse.

12 Abiinformatsiooni kasutamine hinnangutes

Eeldame, et on moodustatud valim vastavalt mingisugusele valikudisainile, on saadud andmed ning ees ootab hindamine. Siiani on meil kasutuses olnud nihketa hinnang kujul

$$\hat{t} = \sum_U \frac{y_i I_i}{EI_i},$$

mille punktihinnang valimi kaudu on järgmine:

$$\hat{t} = \sum_s w_i y_i, \quad w_i = \frac{k_i}{EI_i}.$$

Kaalud w_i selles hinnangus põhinevad pöördväärtusel EI_i , ehk sõltuvad ainult disainist $I \sim p(k)$.

Osutub, et hinnangut summale $t = \sum_U y_i$ on võimalik muuta täpsemaks varieeruvuse mõttes, kui kasutada abiinformatsiooni ja seda just kaalude moodustamise etapil.

Abiinformatsiooniks loetakse

- tunnuseid, mille väärtused on teada kõikide objektide jaoks üldkogumist;

- abitunnuste summasid (näiteks osakogumite kaupa või lihtsalt terves ÜK-s, näiteks meeste arv);
- osakogumite suuruseid üldkogumis (näiteks kihtide mahud).

12.1 Regressioonimudel üldkogumi jaoks

Olgu y_i uuritava tunnuse väärtus objekti i jaoks, $i \in U$. Ja olgu $x_i = (x_{1i}, \dots, x_{ji})^T$ on abitunnuste vektor, mis on teada iga objekti $i, i \in U$ jaoks.

Eeldame järgmist mudelit üldkogumis:

1. väärtus $y_i, i \in U$ on juhusliku suuruse $Y_i \sim \xi$ realisatsioon (jaotusega ξ);
2. jaotuse ξ momendid on järgmiselt defineeritud:
 - $E_\xi Y_i = x_i^T \beta = \sum_{j=1}^J \beta_j x_{ji}$,
 - $V_\xi Y_i = \sigma_i^2$;
3. x_i mittejuhuslikud.

Mudel ütleb seda, et võrdsete x_i korral üldkogumis väärtus y_i võib varieeruda, kuid see varieeruvus toimub tema keskväertuse $x_i^T \beta$ ümbruses (regressioonijoon) dispersiooniga σ_i^2 .

Antud juhul on $\beta = (\beta_1, \dots, \beta_J)^T$ regressioonikordajate vektor.

Märkame, et regressioonimudel on eeldatud üldkogumi väärtustele $y_i, i \in U$. Kui kõik väärtused oleksid teada, siis saaks regressioonikordajate leidmiseks kasutada tavalist kaalutud vähimruutude hinnangut kujul

$$\hat{\beta} \stackrel{\text{tähistame}}{\equiv} B = \left[\sum_U \frac{x_i x_i^T}{\sigma_i^2} \right]^{-1} \sum_U \frac{x_i y_i}{\sigma_i^2}. \quad (70)$$

Prognoositud väärtused y_i -le on $x_i^T B$ ja jäägid üldkogumi mudeli järgi on

$$E_i = y_i - x_i^T B, i \in U. \quad (71)$$

Märkame, et suurused E_i ja B sõltuvad ÜK väärtustest $y_i, i \in U$ ja seega tundmatud. Neid peab hindama valimist. Paneme samuti tähele, et suurus B koosneb kahe ÜK summa korrutisest:

$$\sum_U \frac{x_i x_i^T}{\sigma_i^2} - \text{maatriksite summa, mõõtmetega } J \times J;$$

$$\sum_U \frac{x_i y_i}{\sigma_i^2} - \text{vektorite summa, mõõtmetega } J \times 1.$$

Neid summasid saame hinnata kasutades ÜHT. Lisaks eeldame, et tegemist on TTA disainiga. Siis saame järgmise hinnangu valimist:

$$\hat{B} = \left[\sum_s \frac{x_i x_i^T}{\sigma_i^2 \pi_i} \right]^{-1} \sum_s \frac{x_i y_i}{\sigma_i^2 \pi_i}. \quad (72)$$

Toetudes sellele hinnangule saab arvutada valimist leitud prognoosihinnanguid:

$$\hat{y}_i = x_i^T \hat{B}, i \in U. \quad (73)$$

Mudeli jääkide hinnangud on sel juhul:

$$e_i = y_i - \hat{y}_i, i \in s. \quad (74)$$

kus e_i on leitavad ainult valimis.

12.2 Regressioonihinnang

Et regressioonihinnangut saada, kirjutame üldkogumi summa t ümber:

$$t = \sum_U y_i = \sum_U \hat{y}_i + \sum_U (y_i - \hat{y}_i), \quad (75)$$

kus \hat{y}_i on teada kõikide $i \in U$ korral, ja y_i - ainult valimis.

Hindame teise liikme avaldises (75) kasutades nihketa hinnangut ÜHT-st. See viib regressioonihinnangu kujul:

$$\hat{t}_r = \sum_U \hat{y}_i + \sum_s \frac{y_i - \hat{y}_i}{\pi_i}. \quad (76)$$

Näeme, et regressioonihinnang põhineb prognooside summal, millele on liidetud mudel jääkidest koosnev nn korrigeerimisliige.

Sageli, praktilistel põhjustel esitatakse regressioonihinnang kaalude ja kaalusid korrigeeriva kordaja abil. Selleks kirjutatakse regressioonihinnang ümber:

$$\hat{t}_r = \sum_s \frac{y_i}{\pi_i} + \sum_U \hat{y}_i - \sum_s \frac{\hat{y}_i}{\pi_i}. \quad (77)$$

Nüüd avaldise (73) abil saame

$$\begin{aligned} \hat{t}_r &= \sum_s \frac{y_i}{\pi_i} + \sum_U x_i^T \hat{B} - \sum_s \frac{x_i^T \hat{B}}{\pi_i} = \\ &= \sum_s \frac{y_i}{\pi_i} + \left(\sum_U x_i^T \right) \hat{B} - \left(\sum_s \frac{x_i^T}{\pi_i} \right) \hat{B} \\ &= \sum_s \frac{y_i}{\pi_i} + \left(\sum_U x_i^T - \sum_s \frac{x_i^T}{\pi_i} \right) \hat{B}. \end{aligned}$$

Nüüd, kasutades avaldist (72) \hat{B} jaoks saame:

$$\begin{aligned} \hat{t}_r &= \sum_s \frac{y_i}{\pi_i} + \underbrace{\left(\sum_U x_i^T - \sum_s \frac{x_i^T}{\pi_i} \right)}_{1 \times J} \underbrace{\left[\sum_s \frac{x_i x_i^T}{\sigma_i^2 \pi_i} \right]^{-1}}_{J \times J} \underbrace{\sum_s \frac{x_i y_i}{\sigma_i^2 \pi_i}}_{J \times 1} \\ &= \sum_s \frac{y_i}{\pi_i} \underbrace{\left[1 + \left(\sum_U x_i^T - \sum_s \frac{x_i^T}{\pi_i} \right) \left[\sum_s \frac{x_i x_i^T}{\sigma_i^2 \pi_i} \right]^{-1} \sum_s \frac{x_i}{\sigma_i^2} \right]}_{g_{is}} \end{aligned}$$

Lõplikult, regressioonihinnangut saab esitada kujul

$$\hat{t}_r = \sum_s w_i g_{is} y_i, \quad (78)$$

kus

w_i – on valikukaal,

$$g_{is} = 1 + \left(\sum_U x_i^T - \sum_s \frac{x_i^T}{\pi_i} \right) \left[\sum_s \frac{x_i x_i^T}{\sigma_i^2 \pi_i} \right]^{-1} \sum_s \frac{x_i}{\sigma_i^2}. \quad (79)$$

Valemist (79) näeme, et kui

$$\sum_U x_i^T \approx \sum_s \frac{x_i^T}{\pi_i},$$

ehk x -summad on ligikaudu võrdsed nende hinnangutega, siis $g_{is} \approx 1$, ja

$$\hat{t}_r \approx \sum_s w_i y_i.$$

Regressioonihinnangu dispersiooni tuletust antud kursuse raames ei vaadelda. Siin toome ainult valemi.

Teoreem (Regressioonihinnang). Regressioonihinnang ÜK summale $t = \sum_U y_i$ on antud valemiga (75) ja alternatiivse valemiga (78), mille ligikaudne dispaersioon on

$$V(\hat{t}_r) = \sum_U \sum_U \Delta_{ij}(w_i E_i)(w_j E_j) \quad (80)$$

ja dispersioonihinnanguga

$$\hat{V}(\hat{t}_r) = \sum_{i,j \in s} \frac{\Delta_{ij}}{\pi_{ij}} (w_i g_{is} e_i)(w_j g_{js} e_j), \quad (81)$$

kus üldkogumi taseme jäägid E_i on defineeritud valemiga (71), valimist arvatavad jäägid e_i valemiga (74) ja g -kaalud valemiga (79).

Märkus 1. Teoreemi avaldisest (80) on näha, et mida väiksemad on jäägid E_i , seda väiksem on $V(\hat{t}_r)$. Jäägid E_i näitavad, kui hästi regressioonimudel sobib andmetega. Järelikult, mida parem on regressioonimudel y ja x vahel, seda täpsem tuleb hinnang \hat{t}_r .

Märkus 2. Valemist (79) näeme, et tunnuste x_i üksikuid väärtused peame teadma ainult valimisse sattunud objektide jaoks, üldkogumi tasemel piisab summast $\sum_U x_i$. Järelikult, registrist saab kasutada abiinfot ka agregeeritud kujul (üldkogumi summa näol).

Märkus 3. Paneme tähele, et regressioonihinnangu valem sisaldab suurust σ_i^2 , mis pole aga teada. Praktikas kasutatakse erijuhte, mis eeldavad spetsiaalseid struktuure dispersiooni σ_i^2 jaoks. Tänu nendele probleemidele saab vältida.

Vaata ka ingl. keelset lisamaterjali!