


Теорија узорака


Јован Самарџић, 13/2019


Професор: Милан Јовановић

 - дефиниције

 - ознаке

 - теореме

 - докази

 - примери

Година курса: 2021/22

Молим да ми све грешке пријавите преко мејла или друштвених мрежа.

Популација и узорак

деф. **Популација** је скуп који се посматра. Састоји се од **јединки** ω .

Обим популације је број јединки популације, у ознаци **N** . Он је познат.

деф. **Обележје** је карактеристика коју придружујемо свакој јединки популације.
(најчешће нумеричка, а ако није - „кодирамо“)

Са x_1, \dots, x_N означавамо реализоване вредности обележја популације.

Циљ нам је да сазнамо како се обележје понаша, тј. занима нас какву расподелу има. Пошто није практично гледати све x_i засебно, уводимо следећи појам:

деф. **Параметар обележја** је реална функција вредности обележја популације.

Пример: 1) **укупна (тотална) сума**: $t := \sum_{i=1}^N x_i$;

2) **средња вредност**: $\bar{x} := \frac{1}{N} \sum_{i=1}^N x_i$;

3) **дисперзија/ варијанса**: $s^2 := \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$;

4) **стандардна девијација**: $s := \sqrt{s^2}$;

5) **пропорција неког својства**: $p := \frac{\text{бр. јед. са тим својством}}{\text{укупан број јединки}} = \frac{A}{N}$.

Како прикупљамо податке?

деф. **Попис** је одређивање реализ. вр. за све јединке популације. (добивамо x_1, \dots, x_n)

У пракси, то је тешко изводљиво. Зато не испитујемо целу популацију.

деф. **Узорак** је било који подскуп популације.

Обим узорка је број јединки узорка, у ознаци n . Он је познат.

Напомена: узорак мора бити репрезентативан

Како се бира узорак?

1. корак: **узорачки оквир** - листа/списак свих јединки популације.

Нпр. држављане поређамо у листу по ЈМБГ.

Ипак, узорачки оквир и циљна популација се често не поклопе (нпр. неко нема ЈМБГ)

2. корак: **план узорковања** - правило по ком се од свих могућих узорака фиксiranог обима n бира један узорак

При томе, јављају се грешке и то два типа:

а) **неузорачке грешке** - не зависе од избора плана узорковања.

Примери: - узорачки оквир и популација се разликују;

- двосмислена питања на упитницима;

- мултипликација (телефонска анкета + више фиксних телефона у истој кући);

- проблем мерења.

б) **узорачке грешке** - настају јер закључујемо на основу узорка, а не целе популације.

Напомена: некад можемо да их контролишемо.

Постоје два типа планова узорковања:

деф. **Вероватносни план узорковања** је онај код кога је позната вероватноћа избора било ког узорка фиксиране величине.

Самим тим, позната је вероватноћа појављивања било које јединке у извученом узорку. Због тога је могуће контролисати узорачку грешку.

Постоје 4 типа: 1) **прост случајан узорак**;

2) **статификован узорак**;

3) **кластер узорак**;

4) **систематски узорак**.

У супротном, имамо **невероватносни план узорковања**

На пример: - узорковање само из дела доступне популације;
- узорковање одокативно
- узорковање само оних који то желе, а не свих.

Пример: Имамо оквир $U = \{1, 2, 3, 4\}$.

Узорци обима 2 су: $S_1 = \{1, 2\}$ $S_4 = \{2, 3\}$
 $S_2 = \{1, 3\}$ $S_5 = \{2, 4\}$
 $S_3 = \{1, 4\}$ $S_6 = \{3, 4\}$

Скуп свих узорака S се састоји од S_1, \dots, S_6 . Нама треба један од њих:

$$\left. \begin{array}{l} P(S_1) = 1/3 \\ P(S_2) = 1/6 \\ P(S_3) = P(S_4) = P(S_5) = 0 \\ P(S_6) = 1/2 \end{array} \right\} S: \begin{pmatrix} S_1 & S_2 & S_3 & S_4 & S_5 & S_6 \\ 1/3 & 1/6 & 0 & 0 & 0 & 1/2 \end{pmatrix}$$
 - ово је пример једног вероватносног плана узорковања

$\pi_i := \sum_{S: i \in S} P(S)$ - вероватноћа да се јединка i нађе у узорку који изаберемо

$$\pi_1 = P(S_1) + P(S_2) + P(S_3) = 1/2 ;$$

$$\pi_2 = P(S_1) + P(S_4) + P(S_5) = 1/3 ;$$

$$\pi_3 = P(S_2) + P(S_4) + P(S_6) = 2/3 ;$$

$$\pi_4 = P(S_3) + P(S_5) + P(S_6) = 1/2 .$$

Оцена и особине оцена.

$$\mathbb{R}^n \rightarrow \mathbb{R}$$

деф. **Статистика** је било која Борелова функција која не зависи од непознатих параметара.

деф. **Оцена** је статистика која процењује параметар.

Пример: Настављамо претх. пример.

За оцену статистике t (укупна сума), предлаже се оцена $\hat{t} = N \cdot \bar{x}$.

Дате су реализ. вр:

| i | 1 | 2 | 3 | 4 |
|-------|---|---|---|---|
| x_i | 2 | 5 | 3 | 6 |

(ми овде знамо да је $t=16$)

| | | | | |
|---------------------------------------|---------------|------------------------------------|---|------------------|
| $\bar{x}_{S_1} = \frac{2+5}{2} = 3.5$ | \Rightarrow | $\hat{t}_{S_1} = 4 \cdot 3.5 = 14$ | и | $P(S_1) = 1/3$; |
| $\bar{x}_{S_2} = 2.5$ | \Rightarrow | $\hat{t}_{S_2} = 10$ | и | $P(S_2) = 1/6$; |
| $\bar{x}_{S_3} = 4$ | \Rightarrow | $\hat{t}_{S_3} = 16$ | и | $P(S_3) = 0$; |
| $\bar{x}_{S_4} = 4$ | \Rightarrow | $\hat{t}_{S_4} = 16$ | и | $P(S_4) = 0$; |
| $\bar{x}_{S_5} = 5.5$ | \Rightarrow | $\hat{t}_{S_5} = 22$ | и | $P(S_5) = 0$; |
| $\bar{x}_{S_6} = 4.5$ | \Rightarrow | $\hat{t}_{S_6} = 18$ | и | $P(S_6) = 1/2$. |

$\hat{t}: \begin{pmatrix} 10 & 14 & 16 & 18 & 22 \\ 1/6 & 1/3 & 0 & 1/2 & 0 \end{pmatrix}$ - ово је узорачка расподела оцено \hat{t} .

Приметимо $E\hat{t} = 15\frac{1}{3} \neq \underline{16}$, тј. $E\hat{t} \neq t$.

деф. Нека је $S: \begin{pmatrix} s_1 & \dots & s_b \\ P(s_1) & \dots & P(s_b) \end{pmatrix}$ један план узорковања.

Оцена $\hat{\theta}$ је **непристрасна оцена** параметра θ ако $E(\hat{\theta}) = \theta$. ($E(\hat{\theta}) = \sum \hat{\theta}(s_i) \cdot P(s_i)$)

Дакле, оцена \hat{t} из примера није непристрасна, тј. **пристрасна** је.

деф. **Пристрасност оцене** је $B(\hat{\theta}) := E(\hat{\theta}) - \theta$.

деф. **Средње квадратна грешка оцене** је $MSE(\hat{\theta}) = E((\hat{\theta} - \theta)^2)$.


$$\begin{aligned} \text{Приметимо: } MSE(\hat{\theta}) &= E((\hat{\theta} - \theta)^2) = E((\hat{\theta} - E\hat{\theta} + E\hat{\theta} - \theta)^2) \\ &= E((\hat{\theta} - E\hat{\theta})^2 + 2(\hat{\theta} - E\hat{\theta})(E\hat{\theta} - \theta) + (E\hat{\theta} - \theta)^2) \\ &= E(\hat{\theta} - E\hat{\theta})^2 + 2(E\hat{\theta} - \theta) E(\hat{\theta} - E\hat{\theta}) + E(E\hat{\theta} - \theta)^2 \\ &= D(\hat{\theta}) + 2(E\hat{\theta} - \theta) (E\hat{\theta} - \frac{E(E\hat{\theta})}{E\hat{\theta}}) + (E\hat{\theta} - \theta)^2 \\ &= D(\hat{\theta}) + (B(\hat{\theta}))^2 \end{aligned}$$

Лема 1: $MSE(\hat{\theta}) = D(\hat{\theta}) + (B(\hat{\theta}))^2$.

Доказ: управо извели

Последица: Ако је $\hat{\theta}$ непристрасна $\Rightarrow MSE(\hat{\theta}) = D(\hat{\theta})$.

Напомена: Од две оцене, боља она са мањим MSE.

A hand-drawn teal diamond pattern, consisting of multiple overlapping diamond shapes, covers the entire page. The pattern is centered and fills most of the space, leaving a clear area in the middle for the text.

I deo

деф. **Прост случајан узорак** има најједноставнији план узорковања.
У овом случају, вероватноће избора било ког узорка обима n су једнаке.

Одавде следи и да је једнака вероватноћа да било која јединка буде у узорку.
Другим речима, сви π_i су једнаки.

Обрнуто не важи: постоје планови узорковања са истим π_i ,
али различитим вероватноћама избора било ког узорка обима n . (пример: стр. 22)

1. Прост случајан узорак без понављања

деф. **Прост случајан узорак без понављања** је онај у коме су све јединке у узорку различите.

Својства: - укупан број узорака: $\binom{N}{n}$ ($\{1,2\}$ и $\{2,1\}$ су исти узорак)

- вероватноћа избора узорка: $\frac{1}{\binom{N}{n}}$

- $\pi_i = \frac{\binom{N-1}{n-1}}{\binom{N}{n}} = \frac{n}{N}$ не зависи од $i \Rightarrow$ сви π_i јесу једнаки.

- $\pi_{ij} = \frac{\binom{N-2}{n-2}}{\binom{N}{n}} = \frac{n(n-1)}{N(N-1)}$ ($\pi_{ij} = P\{\text{појављивање } i\text{-те и } j\text{-те јединке у узорку}\}$)

(*) $x_1, \dots, x_i, \dots, x_n$
с одмах га
вирамо

Како извлачимо овај узорак?

I начин: одједном од N извучемо n : $\{x_1, \dots, x_n\}$;

II начин: један по један, са избацивањем.

← коректност

$$P(\{x_{i_1}, \dots, x_{i_n}\}) = \frac{n}{N} \cdot \frac{n-1}{N-1} \cdot \dots \cdot \frac{1}{N-n+1} = \frac{n!}{\frac{N!}{(N-n)!}} = \frac{1}{\binom{N}{n}}$$

У пракси, бирање ПСУ се врши помоћу **таблица случајних бројева**.

Оне су већ генерисане (рачунарски) и прошле су све тестове случајности.

Теорема 1: 1) За ПСУ без понављања, оцена $\bar{x}_n = \frac{1}{n} \sum_{i \in S} x_i$ је непристрасна оцена за \bar{x} ; (S-узорак)

2) За дисперзију те оцене важи: $D(\bar{x}_n) = \frac{s^2}{n} \left(1 - \frac{n}{N}\right)$;

3) Непристрасна оцена те дисперзије је $\widehat{D}(\bar{x}_n) = \frac{s_n^2}{n} \left(1 - \frac{n}{N}\right)$, где $s_n^2 = \frac{1}{n-1} \sum_{i \in S} (x_i - \bar{x}_n)^2$.

Доказ: 1) $\bar{x}_n = \frac{1}{n} \sum_{i \in S} x_i = \frac{1}{n} \sum_{i=1}^N x_i \cdot I_i$, при чему су индикатори зависни (ако их је $n > 1$) (остали су 0)

$$I_i = \begin{pmatrix} 0 & 1 \\ 1 - \pi_i & \pi_i \end{pmatrix}, \quad \text{где } \pi_i = \frac{n}{N}$$

$$E(\bar{x}_n) = E\left(\frac{1}{n} \sum_{i=1}^N x_i \cdot I_i\right) = \frac{1}{n} \sum_{i=1}^N x_i \cdot E(I_i) = \frac{1}{n} \sum_{i=1}^N x_i \cdot \frac{n}{N} = \frac{1}{N} \sum_{i=1}^N x_i = \bar{x} \Rightarrow \text{јесте непристрасна};$$

$$2) D(\bar{x}_n) = D\left(\frac{1}{n} \sum_{i=1}^N x_i \cdot I_i\right) = E\left(\left(\frac{1}{n} \sum_{i=1}^N x_i \cdot I_i\right)^2\right) - \left(E\left(\frac{1}{n} \sum_{i=1}^N x_i \cdot I_i\right)\right)^2 \stackrel{1)}{=} \frac{1}{n^2} E\left(\sum_{i=1}^N x_i I_i\right)^2 - \bar{x}^2$$

$$= \frac{1}{n^2} E\left((x_1 I_1 + \dots + x_N I_N)(x_1 I_1 + \dots + x_N I_N)\right) - \bar{x}^2$$

$$(\Delta) (a_1 + \dots + a_n)^2 = \sum_{i=1}^n a_i^2 + \sum_{i \neq j}^n 2 a_i a_j$$

$$\stackrel{(\Delta)}{=} \frac{1}{n^2} E\left(\sum_{i=1}^N x_i^2 I_i^2 + \sum_{i \neq j}^N x_i x_j I_i I_j\right) - \bar{x}^2$$

$$(**) 1) I_i^2 = I_i$$

$$2) I_i I_j = \begin{pmatrix} 0 & 1 \\ 1 - \pi & \pi \end{pmatrix}$$

$$= \frac{1}{n^2} \left(\sum_{i=1}^N x_i^2 \cdot E(I_i^2) + \sum_{i \neq j}^N x_i x_j E(I_i I_j) \right) - \bar{x}^2$$

$$\stackrel{(**)}{=} \frac{1}{n^2} \left(\sum_{i=1}^N x_i^2 \cdot \frac{n}{N} + \sum_{i \neq j}^N x_i x_j \cdot \frac{n(n-1)}{N(N-1)} \right) - \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2$$

$$= \frac{1}{nN} \left(\sum_{i=1}^N x_i^2 + \frac{n-1}{N-1} \sum_{i \neq j}^N x_i x_j \right) - \frac{1}{N^2} \left(\sum_{i=1}^N x_i^2 + \sum_{i \neq j}^N x_i x_j \right)$$

$$\stackrel{\text{изваљачио заједнички}}{\Rightarrow} = \frac{1}{nN^2} \left[(N-n) \sum_{i=1}^N x_i^2 + \left(\frac{N(n-1)}{N-1} - n \right) \sum_{i \neq j}^N x_i x_j \right]$$

$$= \frac{1}{nN^2} \left[(N-n) \sum_{i=1}^N x_i^2 - (N-n) \frac{1}{N-1} \sum_{i \neq j}^N x_i x_j \right] = \frac{N-n}{nN^2} \left[\sum_{i=1}^N x_i^2 - \frac{1}{N-1} \sum_{i \neq j}^N x_i x_j \right]$$

$$\stackrel{(\ominus)}{=} \frac{N-n}{nN^2} \left[\sum_{i=1}^N x_i^2 - \frac{1}{N-1} \sum_{i \neq j}^N x_i x_j - \frac{1}{N-1} \sum_{i=1}^N x_i^2 + \frac{1}{N-1} \sum_{i=1}^N x_i^2 \right]$$

$$= \frac{1}{nN} \left(1 - \frac{n}{N}\right) \left[\frac{1}{N-1} \sum_{i=1}^N x_i^2 - \frac{1}{N-1} \left(\sum_{i=1}^N x_i^2 + \sum_{i \neq j}^N x_i x_j \right) \right]$$

$$= \frac{1}{nN} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i \right)^2 \right] = \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[\sum_{i=1}^N x_i^2 - \frac{1}{N} \left(\sum_{i=1}^N x_i \right)^2 \right]$$

$$= \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[\sum_{i=1}^N x_i^2 - N \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2 \right] = \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[\sum_{i=1}^N x_i^2 - N \bar{x}^2 \right]$$

$$(***) N \bar{x} = \sum_{i=1}^N x_i$$

Битан трик (користи се и касније)

$$(N-1) s^2 = \left(\sum_{i=1}^N x_i^2 - N \bar{x}^2 \right)$$

$$= \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[\sum_{i=1}^N x_i^2 - 2 N \bar{x}^2 + \underline{N \bar{x}^2} \right] \stackrel{(***)}{=} \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[\sum_{i=1}^N x_i^2 - 2 \bar{x} \sum_{i=1}^N x_i + \sum_{i=1}^N \bar{x}^2 \right]$$

N истих сабирака
 $N \bar{x}^2 = \bar{x}^2 + \dots + \bar{x}^2$
N

$$= \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \sum_{i=1}^N [x_i^2 - 2 x_i \bar{x} + \bar{x}^2]$$

$$= \frac{1}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$$

$$= \frac{s^2}{n} \left(1 - \frac{n}{N}\right)$$

(σ^2 је параметар популације, па је у пракси непознат, па зато га оцењујемо)

$$3) E(D(\bar{x}_n)) = E\left(\frac{s_n^2}{n} \left(1 - \frac{n}{N}\right)\right) = \frac{1}{n} \left(1 - \frac{n}{N}\right) E(s_n^2) \stackrel{M}{=} \frac{s^2}{n} \left(1 - \frac{n}{N}\right) = D(\bar{x}_n) \Rightarrow \text{јесте непристрасна};$$

Овде смо користили непристрасност оцене s_n^2 , па то морамо да докажемо:

Лема 1: За псу без понављања, s_n^2 је непристрасна оцена за σ^2 .

Доказ: Показујемо $E(s_n^2) = \sigma^2$.

$$\begin{aligned} E((n-1)s_n^2) &= E\left(\sum_{i \in S} (x_i - \bar{x}_n)^2\right) \stackrel{\text{Ф.Б. бинома}}{=} E\left(\sum_{i \in S} (x_i^2 - 2x_i \bar{x}_n + \bar{x}_n^2)\right) \\ &= E\left(\sum_{i \in S} x_i^2 - 2\bar{x}_n \sum_{i \in S} x_i + \sum_{i \in S} \bar{x}_n^2\right) \\ &= E\left(\sum_{i \in S} x_i^2 - 2n \frac{1}{n} \bar{x}_n \sum_{i \in S} x_i + n \cdot \bar{x}_n^2\right) \quad \leftarrow n \text{ сабирака (исто као пре)} \\ &= E\left(\sum_{i \in S} x_i^2 - 2n \bar{x}_n^2 + n \bar{x}_n^2\right) = E\left(\sum_{i \in S} x_i^2 - n \bar{x}_n^2\right) \\ &= E\left(\sum_{i \in S} x_i^2 I_i - n \bar{x}_n^2\right) = E\left(\sum_{i \in S} x_i^2 I_i\right) - n E(\bar{x}_n^2) \\ &= \sum_{i \in S} x_i^2 E(I_i) - n E(\bar{x}_n^2) \stackrel{E(x^2) = D(x) + (E(x))^2}{=} \sum_{i \in S} x_i^2 \frac{n}{N} - n(D(\bar{x}_n) + (E(\bar{x}_n))^2) \\ &= \frac{n}{N} \sum_{i \in S} x_i^2 - n \left(\frac{s^2}{n} \left(1 - \frac{n}{N}\right) + \bar{x}^2\right) = \frac{n}{N} \sum_{i \in S} x_i^2 - n \bar{x}^2 - s^2 \left(1 - \frac{n}{N}\right) \\ &= \frac{n}{N} \left(\sum_{i \in S} x_i^2 - N \bar{x}^2\right) - s^2 \left(1 - \frac{n}{N}\right) = \frac{n}{N} (N-1) s^2 - s^2 \left(1 - \frac{n}{N}\right) \\ &= s^2 (n-1) \end{aligned}$$

$N\bar{x} = \sum x_i$
и пишемо $-2 \cdot 1$
па наместимо на s^2
(као у ***) чназад

$$\text{Дакле: } E((n-1)s_n^2) = s^2(n-1) \Rightarrow E(s_n^2) = s^2.$$

деф. **Стандардна грешка** оцене $\hat{\theta}$ је $SE(\hat{\theta}) := \sqrt{MSE(\hat{\theta})}$

Последица: Стандардна грешка оцене \bar{x}_n је $\frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}$.

Напомена: За оцену стандардне грешке, узима се $\frac{s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}$.

Ова оцена није непристрасна! (Без обзира што је $E(s_n^2) = \sigma^2$, не важи $E(s_n) = \sigma$)

Теорема 2: 1) За ПСУ без понављања, оцена $\hat{t} = N \cdot \bar{x}_n$ је непристрасна оцена за t ;

2) За дисперзију те оцене важи: $D(\hat{t}) = N^2 \cdot \frac{S^2}{n} \left(1 - \frac{n}{N}\right)$;

3) Непристрасна оцена те дисперзије је $\hat{D}(\hat{t}) = N^2 \cdot \frac{s_n^2}{n} \left(1 - \frac{n}{N}\right)$.

Доказ: тривијално следи из Т1 и чињенице $t = N \cdot \bar{x}$.

2.

Оцена пропорције на основу ПСУ без понављања

Подсетимо се из [0]: $p = \frac{A}{N}$.

деф. **Узорачка пропорција** је $p_n = \frac{a}{n}$, a - бр. јединки узорка са датим својством.

Кодирање: Ово важи универзално! (без и са понављањем).

$$x_i = \begin{cases} 1, & \text{има то својство} \\ 0, & \text{нема} \end{cases} \quad - \text{ обележје које користимо} \quad \Rightarrow \quad A = \sum x_i$$

$$* \quad p = \frac{\sum x_i}{N} = \bar{x}$$

$$* \quad s^2 = \frac{1}{N-1} \sum (x_i - \bar{x})^2 = \frac{1}{N-1} \left(\sum x_i^2 - N\bar{x}^2 \right)$$

$$\begin{aligned} x_i^2 = x_i &\approx \frac{1}{N-1} \left(\sum x_i - N\bar{x}^2 \right) = \frac{1}{N-1} (Np - Np^2) \\ &= \frac{N}{N-1} p(1-p) \end{aligned}$$

Теорема 1: 1) За ПСУ без понављања, оцена $p_n = \frac{a}{n}$ је непристрасна оцена за p ;

2) За дисперзију те оцене важи: $D(p_n) = \frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}$;

3) Непристрасна оцена те дисперзије је $\widehat{D}(p_n) = \frac{N-n}{N} \cdot \frac{p_n(1-p_n)}{n-1}$.

Доказ: 1) По ПТ.1: $E(\bar{x}_n) = \bar{x} \stackrel{\text{коп.}}{\Rightarrow} E(p_n) = p \Rightarrow$ јесте непристрасна;

2) $D(p_n) = \frac{s^2}{n} \left(1 - \frac{n}{N}\right) \stackrel{\text{коп.}}{=} \frac{1}{n} \cdot \frac{N}{N-1} p(1-p) \left(1 - \frac{n}{N}\right) = \frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}$;

3) $\widehat{D}(p_n) = \frac{s_n^2}{n} \left(1 - \frac{n}{N}\right) \stackrel{\text{коп.}}{=} \frac{n}{n-1} \cdot \frac{p_n(1-p_n)}{n} \cdot \frac{N-n}{n} = \frac{N-n}{N} \cdot \frac{p_n(1-p_n)}{n-1} \Rightarrow$ јесте непристрасна.

Лема 1: $s_n^2 = \frac{n}{n-1} p_n(1-p_n)$

Доказ: $s_n^2 = \frac{1}{n-1} \sum_{i \in S} (x_i - \bar{x}_n)^2 = \frac{1}{n-1} (\sum_{i \in S} x_i^2 - n\bar{x}_n^2)$
 $= \frac{1}{n-1} (\sum_{i \in S} x_i - n\bar{x}_n^2) = \frac{1}{n-1} (np_n - np_n^2)$
 $= \frac{n}{n-1} p_n(1-p_n).$

Теорема 2: 1) За ПСУ без понављања, оцена $\hat{A} = N \cdot p_n$ је непристрасна оцена за A ;

2) За дисперзију те оцене важи: $D(\hat{A}) = N^2 \cdot \frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}$;

3) Непристрасна оцена те дисперзије је $\widehat{D}(\hat{A}) = N \cdot (N-n) \cdot \frac{p_n(1-p_n)}{n-1}$.

Доказ: тривијално следи из ПТ и чињенице $A = N \cdot p$.

3. Прост случајан узорак са понављањем

деф. **Прост случајан узорак са понављањем** је онај у коме се јединке у узорку могу понављати.

Својства: - укупан број узорака: N^n

- вероватноћа избора узorca: $\frac{1}{N^n}$ ($= \frac{1}{N} \cdot \dots \cdot \frac{1}{N}$)

Напомена: Редослед је битан! (нпр. (1,3,4,5) и (5,1,3,4) су различити узорци)
Зато кажемо да су ово уређени узорци.

Такође, нпр. $P(\{1,1,1,1\}) \neq P(\{1,3,4,5\})$.
↳ само (1,1,1,1) ↳ може бити (1,3,4,5), (5,1,3,4)...

Теорема 1: 1) За ПСУ са понављањем, оцена $\bar{x}_n = \frac{1}{n} \sum x_i$ је непристрасна оцена за \bar{x} ;

2) За дисперзију те оцене важи: $D(\bar{x}_n) = \frac{s^2}{n} (1 - \frac{1}{N})$;

3) Непристрасна оцена те дисперзије је $\hat{D}(\bar{x}_n) = \frac{s_n^2}{n}$.

Доказ: Имамо узорак обима n : (x_1, \dots, x_n) - x_i су независне и имају исту расподелу као обележје X .

$$X: \begin{pmatrix} x_1 & \dots & x_N \\ 1/N & \dots & 1/N \end{pmatrix}$$

случајно \Rightarrow све једнако $\Rightarrow \frac{1}{N}$

Рачунамо: * $EX = \frac{1}{N} \cdot \sum x_i = \bar{x}$;

$$* DX = EX^2 - (EX)^2 = \frac{1}{N} (x_1^2 + \dots + x_n^2) - \bar{x}^2$$

$$= \frac{1}{N} [\sum x_i^2 - N\bar{x}^2] = \frac{N-1}{N} \cdot s^2. \quad \text{окуп: } (N-1)s^2 = (\sum x_i^2 - N\bar{x}^2)$$

1) $E(\bar{x}_n) = E(\frac{1}{n} \sum x_i) = \frac{1}{n} \cdot \sum EX_i = \frac{1}{n} \sum EX = \frac{1}{n} \cdot n \cdot EX = \bar{x} \Rightarrow$ јесте непристрасна;

2) $D(\bar{x}_n) = D(\frac{1}{n} \sum x_i) \stackrel{\text{нес.}}{=} \frac{1}{n^2} \sum D(x_i) = \frac{n}{n^2} DX = \frac{s^2}{n} (1 - \frac{1}{N})$;

$$3) E(D(\hat{\bar{x}}_n)) = E\left(\frac{s_n^2}{n}\right) = \frac{1}{n} E(s_n^2) \stackrel{M1}{=} \frac{1}{n} \cdot \frac{N-1}{N} S^2 = D(\bar{X}_n) \Rightarrow \text{јесте непристрасна.}$$

Лема 1: За ПСУ са понављањем: $E(s_n^2) = \frac{N-1}{N} \cdot S^2$

$$\begin{aligned} \text{Доказ: } E(s_n^2) &= E\left(\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)^2\right) = \frac{1}{n-1} E\left(\sum_{i=1}^n x_i^2 - n \cdot \bar{x}_n^2\right) = \frac{1}{n-1} (nE(x^2) - nE(\bar{x}_n^2)) \\ &= \frac{1}{n-1} \left[n \cdot (DX + (EX)^2) - n \cdot (D\bar{X}_n + (E\bar{X}_n)^2) \right] \\ &\stackrel{1,2)}{=} \frac{n}{n-1} \left[\frac{N-1}{N} S^2 + \bar{x}^2 - \frac{S^2}{n} \left(1 - \frac{1}{N}\right) - \bar{x}^2 \right] \\ &= \frac{N-1}{N} \cdot \frac{n}{n-1} \cdot S^2 \left(1 - \frac{1}{N}\right) = \frac{N-1}{N} \cdot S^2. \end{aligned}$$

Теорема 2: 1) За ПСУ са понављањем, оцена $\hat{t} = N \cdot \bar{x}_n$ је непристрасна оцена за t ;

2) За дисперзију те оцене важи: $D(\hat{t}) = \frac{N(N-1)}{n} \cdot S^2$;

3) Непристрасна оцена те дисперзије је $\hat{D}(\hat{t}) = \frac{N^2 s_n^2}{n}$.

Доказ: тривијално следи из Т1 и чињенице $\hat{t} = N \cdot \bar{x}$.

* Упоредимо оцене за случајеве ПСУ без и са понављањем:

Знамо да је боља оцена која има мање MSE.

Како су наше оцене непристрасне, гледамо која има мању дисперзију.

$$D_{\text{вр}}(\bar{x}_n) = \frac{N-n}{N-1} D_{\text{сп}}(\bar{x}_n) \Rightarrow D_{\text{вр}}(\bar{x}_n) < D_{\text{сп}}(\bar{x}_n);$$

Закључак: Оцене без понављања су боље од оцена са понављањем код коначних популација.
 $N < \infty$

4.

Оцена пропорције на основу ПСУ са понављањем

Кодирање: исто као у [2].

Теорема 1: 1) За ПСУ са понављањем, оцена $p_n = \frac{a}{n}$ је непристрасна оцена за p ;

2) За дисперзију те оцене важи: $D(p_n) = \frac{p(1-p)}{n}$;

3) Непристрасна оцена те дисперзије је $\hat{D}(p_n) = \frac{p_n(1-p_n)}{n-1}$.

Доказ: 1) $E(\bar{x}_n) = \bar{x} \stackrel{\text{код}}{\Rightarrow} E(p_n) = p$;

2) $D(p_n) = \frac{s^2}{n} \left(1 - \frac{1}{N}\right) \stackrel{\text{код}}{=} \frac{N}{N-1} \cdot \frac{p(1-p)}{n} \cdot \frac{N-1}{N} = \frac{p(1-p)}{n}$;

3) $\hat{D}(p_n) = \frac{s_n^2}{n} \stackrel{[2].n}{=} \frac{n}{n-1} \cdot \frac{p_n(1-p_n)}{n} = \frac{p_n(1-p_n)}{n-1}$.

Теорема 2: 1) За ПСУ са понављањем, оцена $\hat{A} = N \cdot p_n$ је непристрасна оцена за A ;

2) За дисперзију те оцене важи: $D(\hat{A}) = N^2 \cdot \frac{p(1-p)}{n}$;

3) Непристрасна оцена те дисперзије је $\hat{D}(\hat{A}) = N^2 \cdot \frac{p_n(1-p_n)}{n-1}$.

Доказ: тривијално следи из т1 и чињенице $A = N \cdot p$.

5.

Интервали поверења на основу ПСУ без понављања

По сада смо гледали искључиво тачкасте оцене, а сад гледамо интервалне.

деф. Са I означавамо интервал поверења, а $1-\alpha$ нам је ниво поверења.

То значи $P\{\theta \in I\} = 1-\alpha$, где је θ непознато.

Интервал поверења за \bar{x}

Хајек је показао (један облик) ЦГТ за ПСУ без понављања:

За довољно велико $n, N, N-n$ код ПСУ без понављања важи:

$n \geq 30$

$$X^* = \frac{\bar{x}_n - \bar{x}}{\frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}} \approx \mathcal{N}(0,1)$$

$\xrightarrow{E\bar{x}_n}$
 $\xrightarrow{\sqrt{D\bar{x}_n}}$

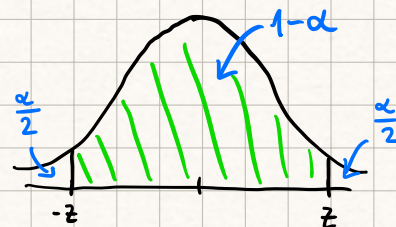
Наравно, постоји ∞ интервала поверења, па морамо неки да изаберемо.

Посматрамо тзв. центрирани интервал поверења:

$$P\{|X^*| \leq z\} = P\{-z \leq X^* \leq z\} = 1-\alpha$$

$$P\left\{-z \leq \frac{\bar{x}_n - \bar{x}}{\frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}} \leq z\right\} = 1-\alpha$$

$$P\left\{\bar{x}_n - z \cdot \frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}} \leq \bar{x} \leq \bar{x}_n + z \cdot \frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}\right\}$$



$$\Rightarrow I_{\bar{x}} = \left[\bar{x}_n - z \cdot \frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}, \bar{x}_n + z \cdot \frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}} \right], \quad \text{где } F_{X^*}(z) = 1 - \frac{\alpha}{2}.$$

Ово је интервал поверења величине $1-\alpha$ за ср. вр. обележја популације.

Проблем: s је параметар популације и непознат је. Зато га оцењујемо са s_n :

$$I_{\bar{x}} = \left[\bar{x}_n - z \cdot \frac{s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}, \bar{x}_n + z \cdot \frac{s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}} \right], \quad \text{где } F_{X^*}(z) = 1 - \frac{\alpha}{2}. \quad (*)$$

Ово је апроксимативни интервал поверења за ср. вр. обележја популације.

За $n < 30$: користимо Студентову t_{n-1} расподелу. ($n-1$ степена слободне)

$$I_{\bar{x}} = \left[\bar{x}_n - t \cdot \frac{s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}, \bar{x}_n + t \cdot \frac{s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}} \right], \quad \text{где } F_{t_{n-1}}(t) = 1 - \frac{\alpha}{2}.$$

Интервал поверења за t

Користимо чињеницу да је $\hat{t} = N \cdot \bar{x}_n$.

Зато интервал (*) множимо са N : $I_{N \cdot \bar{x}} = \left[N \cdot \bar{x}_n - z \cdot \frac{N \cdot s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}, N \cdot \bar{x}_n + z \cdot \frac{N \cdot s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}} \right]$,

$$\Rightarrow I_t = \left[\hat{t} - z \cdot \frac{N \cdot s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}, \hat{t} + z \cdot \frac{N \cdot s_n}{\sqrt{n}} \sqrt{1 - \frac{n}{N}} \right]$$

При томе, за $n \geq 30$: $F_{\chi^2}(z) = 1 - \frac{\alpha}{2}$;

за $n < 30$: $F_{t_{n-1}}(t) = 1 - \frac{\alpha}{2}$.

Напомена: исти резултат бисмо добили истим поступком као за \bar{x} , само крећемо од $\frac{\hat{t} - t}{\frac{s}{\sqrt{n}} \sqrt{N(N-n)}} \approx N(0,1)$.

Интервал поверења за p

Аналогно почетном поступку, добија се:

$$I_p = \left[p_n - z \cdot \sqrt{\frac{p_n(1-p_n)}{n-1} \left(1 - \frac{n}{N}\right)} - \frac{1}{2n}, p_n + z \cdot \sqrt{\frac{p_n(1-p_n)}{n-1} \left(1 - \frac{n}{N}\right)} + \frac{1}{2n} \right].$$

✿ - корекција (x_n - непрекидно, p_n дискретно)

Множењем I_p са N добијамо интервал поверења за A :

$$I_A = \left[\hat{A} - z \cdot \sqrt{\frac{p_n(1-p_n)}{n-1} (N-n)} - \frac{N}{2n}, \hat{A} + z \cdot \sqrt{\frac{p_n(1-p_n)}{n-1} (N-n)} + \frac{N}{2n} \right].$$

6.

Одређивање обима ПСУ без понављања

Ако имамо већи обим, имамо и прецизније резултате.

Са друге стране, јављају се веће неузорачке грешке, као и трошкови.

Наш циљ је да одредимо оптимално n .

Уведимо ознаке: θ : параметар који оцењујемо;

$\hat{\theta}$: оцена параметра;

d : "подношљиво" одступање ($d > 0$);

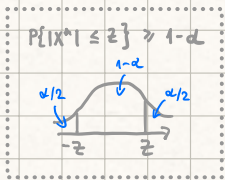
α : ризик да се премаши d (обично мало: 0.05, 0.01, 0.1).

} унапред задати

Дакле: $P\{|\hat{\theta} - \theta| > d\} \leq \alpha$

Одређивање обима за оцењивање \bar{X}

$$P\{|\bar{X}_n - \bar{X}| > d\} \leq \alpha \Rightarrow P\{|\bar{X}_n - \bar{X}| \leq d\} \geq 1 - \alpha \Rightarrow P\left\{\underbrace{\frac{|\bar{X}_n - \bar{X}|}{\frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}}}_{\chi^* \sim N(0,1)} \leq \frac{d}{\frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}}\right\} \geq 1 - \alpha$$



$$\Rightarrow \frac{d}{\frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}} \geq z, \quad \text{где } \Phi(z) = 1 - \frac{\alpha}{2} \Rightarrow \frac{d}{z} \geq \frac{s}{\sqrt{n}} \sqrt{1 - \frac{n}{N}}$$

$$\Rightarrow \frac{d^2}{z^2} \geq \frac{s^2}{n} \left(1 - \frac{n}{N}\right)$$

$$\Rightarrow \frac{d^2}{z^2} \geq \frac{s^2}{n} - \frac{s^2}{N} \Rightarrow \frac{s^2}{n} \leq \frac{d^2}{z^2} + \frac{s^2}{N}$$

Теорема 1: $n \geq \frac{1}{\frac{d^2}{s^2 z^2} + \frac{1}{N}}$, где $\Phi(z) = 1 - \frac{\alpha}{2}$

Напомена: d, α, N су познати, z прочитамо из таблице.

Проблем: s^2 је непознато.

Не можемо га оценити са s_n^2 , јер узорак још немамо.

Решење: Узимамо узорак из сличних истраживања из прошлости или направимо прелиминарно истраживање.

Тако добијамо s_n^2 .

Одређивање обима за оцењивање t

$$P\{|\hat{t}-t| > d\} \leq \alpha \Rightarrow P\{|\hat{t}-t| \leq d\} \geq 1-\alpha \Rightarrow P\left\{\frac{|\hat{t}-t|}{\frac{\Delta N}{\sqrt{n}} \sqrt{1-\frac{n}{N}}} \leq \frac{d}{\frac{\Delta N}{\sqrt{n}} \sqrt{1-\frac{n}{N}}}\right\} \geq 1-\alpha$$

$\sim N(0,1)$

$$\Rightarrow \frac{d}{\frac{\Delta N}{\sqrt{n}} \sqrt{1-\frac{n}{N}}} \geq z, \quad \text{где } \Phi(z) = 1 - \frac{\alpha}{2} \Rightarrow \frac{d}{z} \geq \frac{\Delta N}{\sqrt{n}} \sqrt{1-\frac{n}{N}}$$

$$\Rightarrow \frac{d^2}{z^2} \geq \frac{N^2 \Delta^2}{n} \cdot \frac{N-n}{N}$$

$$\Rightarrow \frac{d^2}{z^2} \geq \frac{N^2 \Delta^2}{n} - N \Delta^2 \Rightarrow \frac{N^2 \Delta^2}{n} \leq \frac{d^2}{z^2} + N \Delta^2$$

Теорема 2: $n \geq \frac{1}{\frac{d^2}{z^2 N^2 \Delta^2} + \frac{1}{N}}$, где $\Phi(z) = 1 - \frac{\alpha}{2}$

Напомена: Проблем непознатог Δ решавамо на исти начин.

Одређивање обима за оцењивање p

$$P\{|p_n - p| > d\} \leq \alpha \Rightarrow P\{|p_n - p| \leq d\} \geq 1-\alpha \Rightarrow P\left\{\frac{|p_n - p|}{\sqrt{\frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}}} \leq \frac{d}{\sqrt{\frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}}}\right\} \geq 1-\alpha$$

$\sim N(0,1)$

$$\Rightarrow \frac{d}{\sqrt{\frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}}} \geq z, \quad \text{где } \Phi(z) = 1 - \frac{\alpha}{2} \Rightarrow \frac{d}{z} \geq \sqrt{\frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}}$$

$$\Rightarrow \frac{d^2}{z^2} \geq \frac{N-n}{N-1} \cdot \frac{p(1-p)}{n}$$

$$\Rightarrow \frac{d^2}{z^2} \geq \frac{N}{N-1} \cdot \frac{p(1-p)}{n} - \frac{p(1-p)}{N-1} \Rightarrow \frac{N p(1-p)}{(N-1)n} \leq \frac{d^2}{z^2} + \frac{p(1-p)}{N-1}$$

Теорема 3: $n \geq \frac{1}{\frac{d^2}{z^2 \cdot \frac{N}{N-1} p(1-p)} + \frac{1}{N}}$, где $\Phi(z) = 1 - \frac{\alpha}{2}$

Напомена: p нам је непознато. Може као и пре, а може и на други начин:

Искористимо да важи $p(1-p) \leq 1/4$.

Мало грубља апроксимација, али олакшава рачун

Тада је: $n \geq \frac{1}{\frac{4d^2}{z^2 \cdot \frac{N}{N-1}} + \frac{1}{N}}$, где $\Phi(z) = 1 - \frac{\alpha}{2}$

Овим смо завршили ПСУ.

Поставља се питање када треба користити ПСУ?

У случајевима када имамо добар узорачки оквир,
али скоро никакве додатне информације по којим бисмо јединке могли да групишемо.

Пример: Имамо списак студената целог универзитета,
али само по имену и презимену.

7.

Методи избора узорка са неједнаким вероватноћама избора јединки

деф. Уводимо помоћно обележје: свакој јединки додељујемо величину M_i . ($i = \overline{1, N}$)

При томе, $\sum_{i=1}^N M_i = M$.

деф. $p_i := \frac{M_i}{M}$ је вероватноћа избора i -те јединке у једном бирању.

Напомена: При томе, очигледно важи $\sum_{i=1}^N p_i = 1$.

а) са понављањем: - узмемо, па вратимо;
- тиме почетна расподела p_i -ева остаје иста ($p_i = \frac{M_i}{M}$)

б) без понављања: - извучена i -та ($p_i = \frac{M_i}{M}$)
- избацимо је
- извлачимо j -ту ($p_j = \frac{M_j}{M - M_i}$)
- избацимо је
:

Наводимо методе избора оваквог узорка:

1) Метод кумуланте:

Из таблице случ. бр.¹⁾ бирамо број $V = \overline{1, M}$.

$V \in [1, M_1] \Rightarrow$ извукли смо јединку 1;

$V \in [M_1+1, M_1+M_2] \Rightarrow$ извукли смо јединку 2;

\vdots

$V \in [M_1+\dots+M_{k-1}+1, M_1+\dots+M_k] \Rightarrow$ извукли смо јединку k ;

\vdots

$V \in [M_1+\dots+M_{n-1}+1, \underbrace{M_1+\dots+M_n}_M] \Rightarrow$ извукли смо јединку n .

Границе се зову **кумуланте**. $(M_1, M_1+M_2, M_1+M_2+M_3, \dots)$

Метод је коректан:
$$p_i = \frac{M_1+\dots+M_i - (M_1+\dots+M_{i-1})}{M} = \frac{M_i}{M}$$

Напомена: Све ово је за узорак са понављањем.

За без понављања: ако извучемо неки од раније - одбацујемо га.

Мана: код великих популација, тешко рачунамо кумуланте.

Зато уводимо други метод.

2) Лахиријев метод:

Означимо са $K = \max M_i$.

Бирамо уређени пар: (i, R) - где $i \in \{1, 2, \dots, N\}$ и $R \in \{1, 2, \dots, K\}$

Ако је $R \leq M_i \Rightarrow i$ -та јединка је ушла у узорак;
 $R > M_i \Rightarrow$ није ушла.

Поступак понављамо док не добијемо n јединки. (различитих, ако је без понављања)

Метод је коректан:

$p_i = \sum_{s=1}^{\infty} p_i(s)$, где је $p_i(s)$ вероватноћа да i -та јединка буде извучена у s -том покушају
($s-1$ ни једна, па s -ти покушај i -ту)
оно старо

Вани: $p_i(1) = \frac{1}{N} \cdot \frac{M_i}{K}$ (\square - морамо извући i -ту јединку од њих N)
 \blacksquare - мора $R \leq M_i$, а $R = 1, K$

$$p_i(2) = \left(1 - \sum_{i=1}^N \frac{1}{N} \cdot \frac{M_i}{K} \right) \cdot \frac{1}{N} \cdot \frac{M_i}{K}$$

\vdots
 q (није извучена ни једна јединка) (i -та покушај)

Уочимо правило: $p_i(s) = q^{s-1} \cdot \frac{1}{N} \cdot \frac{M_i}{K}$.

Убацимо то: $p_i = \sum_{s=1}^{\infty} \left(q^{s-1} \cdot \frac{1}{N} \cdot \frac{M_i}{K} \right) \stackrel{\sum_{s=1}^{\infty} q^{s-1} = \frac{1}{1-q}}{=} \frac{1}{N} \cdot \frac{M_i}{K} \cdot \frac{1}{1 - \left(1 - \sum_{i=1}^N \frac{1}{N} \cdot \frac{M_i}{K} \right)}$

$$= \frac{1}{N} \cdot \frac{M_i}{K} \cdot \frac{1}{\frac{1}{NK} \sum M_i}$$
$$= \frac{M_i}{M}$$

8.

Hansen - Hurwitz - ове оцене

Гледамо како изгледају оцене код јединке са неједнаким вероватноћама избора, са понављањем.
 Нове оцене ће у индексу имати ознаку **НН**.

Теорема 1: Нека је p_i вероватноћа избора i -те јединке у извлачењу код узорковања са понављањем.

- 1) Оцена $\hat{t}_{НН} = \frac{1}{n} \sum_{i \in S} \frac{x_i}{p_i}$ је непристрасна оцена за t ;
- 2) За дисперзију те оцене важи: $D(\hat{t}_{НН}) = \frac{1}{n} \sum_{i=1}^N p_i \left(\frac{x_i}{p_i} - t \right)^2$;
- 3) Непристрасна оцена те дисперзије је $\widehat{D(\hat{t}_{НН})} = \frac{1}{n(n-1)} \sum_{i \in S} \left(\frac{x_i}{p_i} - \hat{t}_{НН} \right)^2$.

Доказ: 1) $E(\hat{t}_{НН}) = E\left(\frac{1}{n} \sum_{i \in S} \frac{x_i}{p_i}\right) = \frac{1}{n} E\left(\sum_{i \in S} \frac{x_i}{p_i}\right) = \frac{1}{n} E(\sum_{i \in S} y_i) \stackrel{\downarrow}{=} \frac{1}{n} \sum_{i \in S} E y_i$
 $\stackrel{(*)}{=} \frac{1}{n} \cdot n \cdot t = t \Rightarrow$ јесте непристрасна

(*) $y: \left(\frac{x_1}{p_1}, \dots, \frac{x_n}{p_n}\right) \Rightarrow E y = \sum x_i = t$
 И имамо да $i \in S$ (узорку) $\Rightarrow n \cdot t$
 (не популацији)

2) $D(\hat{t}_{НН}) = D\left(\frac{1}{n} \sum_{i \in S} \frac{x_i}{p_i}\right) \stackrel{i.i.d.}{=} \frac{1}{n^2} \sum_{i \in S} D y_i = \frac{1}{n^2} \cdot n \cdot D y = \frac{1}{n} D y$
 $\stackrel{(**) I}{=} \frac{1}{n} \sum \left(\frac{x_i}{p_i} - t\right)^2 \cdot p_i$

(**) Рачунамо на два начина:
 I) $D y = E(y - E y)^2 = \sum \left(\frac{x_i}{p_i} - t\right)^2 \cdot p_i$
 II) $D y = E y^2 - (E y)^2 = \sum \frac{x_i^2}{p_i} - t^2$

3) $E(\widehat{D(\hat{t}_{НН})}) = E\left[\frac{1}{n(n-1)} \sum_{i \in S} \left(\frac{x_i}{p_i} - \hat{t}_{НН}\right)^2\right] = \frac{1}{n(n-1)} E\left[\sum_{i \in S} \left(\frac{x_i}{p_i} - \frac{1}{n} \sum_{i \in S} \frac{x_i}{p_i}\right)^2\right]$
 $= \frac{1}{n(n-1)} E\left[\sum_{i \in S} \left(\frac{x_i^2}{p_i^2} - 2 \frac{x_i}{p_i} \frac{1}{n} \sum_{i \in S} \frac{x_i}{p_i} + \frac{1}{n^2} \left(\sum_{i \in S} \frac{x_i}{p_i}\right)^2\right)\right]$
 $= \frac{1}{n(n-1)} E\left[\sum_{i \in S} \frac{x_i^2}{p_i^2} - \frac{2}{n} \left(\sum_{i \in S} \frac{x_i}{p_i}\right)^2 + \frac{n}{n^2} \left(\sum_{i \in S} \frac{x_i}{p_i}\right)^2\right] = \frac{1}{n(n-1)} E\left[\sum_{i \in S} \frac{x_i^2}{p_i^2} - \frac{1}{n} \left(\sum_{i \in S} \frac{x_i}{p_i}\right)^2\right]$
 $= \frac{1}{n(n-1)} \left[\sum_{i \in S} E\left(\frac{x_i^2}{p_i^2}\right) - \frac{1}{n} E\left(\sum_{i \in S} \frac{x_i}{p_i}\right)^2\right] = \frac{1}{n(n-1)} \left[\sum_{i \in S} E(y_i^2) - \frac{1}{n} E\left(\sum_{i \in S} y_i\right)^2\right]$
 $\stackrel{E^2 = D X + (E X)^2}{=} \frac{1}{n(n-1)} \left[n \cdot \sum_{i \in S} \frac{x_i^2}{p_i^2} \cdot p_i - \frac{1}{n} (D(\sum y_i) + E(\sum y_i)^2)\right]$
 $\stackrel{i.i.d.}{=} \frac{1}{n(n-1)} \left[n \cdot \sum_{i \in S} \frac{x_i^2}{p_i^2} - \frac{1}{n} (n D y + (n E y)^2)\right]$
 $\stackrel{(*)}{=} \frac{1}{n(n-1)} \left[n \cdot \sum_{i \in S} \frac{x_i^2}{p_i^2} - n D y - n t^2\right]$
 $= \frac{1}{n(n-1)} \left[n \cdot \left(\sum_{i \in S} \frac{x_i^2}{p_i^2} - t^2\right) - n D y\right] \stackrel{(**) II}{=} \frac{1}{n(n-1)} [n \cdot D y - n D y]$
 $= \frac{1}{n} \cdot D y \stackrel{(**) I}{=} D(\hat{t}_{НН}) \Rightarrow$ јесте непристрасна

Теорема 2: Нека је p_i вероватноћа избора i -те јединке у извлачењу код узорковања са понављањем.

- 1) Оцена $\bar{x}_{HH} = \frac{1}{n} \sum_{i \in S} \frac{x_i}{p_i}$ је непристрасна оцена за t ;
- 2) За дисперзију те оцене важи: $D(\bar{x}_{HH}) = \frac{1}{n} \sum_{i=1}^N p_i \left(\frac{x_i}{p_i} - \bar{x} \right)^2$;
- 3) Непристрасна оцена те дисперзије је $\hat{D}(\bar{x}_{HH}) = \frac{1}{n(n-1)} \sum_{i \in S} \left(\frac{x_i}{p_i} - \bar{x}_{HH} \right)^2$.

Доказ: тривијално следи из Т1 и чињенице $\bar{x} = \frac{1}{N} \cdot t$.

Напомена: Оцене за ПСУ са понављањем су спец. случај НН.

Доказ: Користимо да је у ПСУ $p_i = \frac{1}{N}$.

$$\bar{x}_{HH} = \frac{1}{N \cdot n} \sum_{i \in S} \frac{x_i}{p_i} = \frac{1}{Nn} \sum_{i \in S} x_i \cdot N = \frac{1}{n} \sum_{i \in S} x_i = \bar{x}_n$$

Остале такође тривијално

Напомена: Предности оцене \hat{t}_{HH} у односу на ПСУ са понављањем:

- није потребно знати N ;
- није потребно знати све p_i , већ само за оне i које су упале у узорак.

Horvitz - Thompson - ове оцене

За неједнаке вероватноће избора, без понављања, оцене нису неке + су компликоване.
 Ми ћемо радити Horvitz - Thompson - ове, које могу да се користе за сваки план узорковања. (и са и без)

Теорема 1: Нека је π_i вероватноћа избора i -те јединке.
 Нека је π_{ij} ($i \neq j$) вероватноћа избора i -те и j -те јединке.
 Такође, нека $\pi_i > 0$, $\pi_{ij} > 0$. ($\forall i, j = \overline{1, N}$, $i \neq j$)

1) Оцена $\hat{t}_{HT} = \sum_{i=1}^{\mathcal{J}} \frac{x_i}{\pi_i}$ је непристрасна оцена за t ;

2) За дисперзију те оцене важи: $D(\hat{t}_{HT}) = \sum_{i=1}^{\mathcal{J}} \frac{1-\pi_i}{\pi_i} \cdot x_i^2 + \sum_{i=1}^{\mathcal{J}} \sum_{j=1, j \neq i}^{\mathcal{J}} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} x_i x_j$;

3) Непристрасна оцена те дисперзије је $D(\hat{t}_{HT}) = \sum_{i=1}^{\mathcal{J}} \frac{1-\pi_i}{\pi_i^2} \cdot x_i^2 + \sum_{i=1}^{\mathcal{J}} \sum_{j=1, j \neq i}^{\mathcal{J}} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{x_i x_j}{\pi_i \pi_j}$;

Напомена: \mathcal{J} је број различитих јединки у узорку.

Под \sum мисли се да узимамо различите, што не мора да значи првих \mathcal{J} !

Другим речима: $\hat{t}_{HT} = \sum_{i=1}^{\mathcal{J}} \frac{x_i}{\pi_i} = \sum_{i=1}^N \frac{x_i}{\pi_i} \cdot I_i$

букв. исто што и $i \in S$
ако се нека понови, због ове суме ће се рачунати само једном
 $I_i = \begin{pmatrix} 0 & 1 \\ 1-\pi_i & \pi_i \end{pmatrix}$, $I_{ij} = \begin{pmatrix} 0 & 1 \\ 1-\pi_{ij} & \pi_{ij} \end{pmatrix}$
 \uparrow i -та јединка ушла у узорак

За узорак без понављања важи $\mathcal{J} = n$.

Напомена: $D(\hat{t}_{HT})$ има проблем: $\pi_{ij} - \pi_i \pi_j$ може бити негативно.

Тада је ова оцена бескорисна и праве се (компликованије) пристрасне модификације.

Исто ће важити у [10] за $D_{SSA}(\hat{t}_{HT})$.

Доказ: 1) $E(\hat{t}_{HT}) = E\left(\sum_{i=1}^N \frac{x_i}{\pi_i} \cdot I_i\right) = \sum_{i=1}^N \frac{x_i}{\pi_i} E I_i = \sum_{i=1}^N \frac{x_i}{\pi_i} \pi_i = \sum x_i = t \Rightarrow$ јесте непристрасна;

$$\begin{aligned} 2) D(\hat{t}_{HT}) &= D\left(\sum_{i=1}^N \frac{x_i}{\pi_i} \cdot I_i\right) = E\left[\left(\sum_{i=1}^N \frac{x_i}{\pi_i} \cdot I_i - E\left(\sum_{i=1}^N \frac{x_i}{\pi_i} \cdot I_i\right)\right)^2\right] \\ &= E\left[\left(\sum_{i=1}^N \frac{x_i}{\pi_i} \cdot I_i - \sum_{i=1}^N \frac{x_i}{\pi_i} E I_i\right)^2\right] = E\left[\left(\sum_{i=1}^N \frac{x_i}{\pi_i} \cdot I_i - \sum_{i=1}^N \frac{x_i}{\pi_i} \pi_i\right)^2\right] \\ &= E\left[\left(\sum_{i=1}^N \frac{x_i}{\pi_i} (I_i - \pi_i)\right)^2\right] \end{aligned}$$

$$(\Delta) (a_1 + \dots + a_n)^2 = \sum_i a_i^2 + \sum_{i \neq j} 2 a_i a_j$$

$$\begin{aligned} &\stackrel{(\Delta)}{=} E\left[\sum_{i=1}^N \frac{x_i^2}{\pi_i^2} (I_i - \pi_i)^2 + \sum_{i=1}^N \sum_{j \neq i}^N \frac{x_i}{\pi_i} \frac{x_j}{\pi_j} (I_i - \pi_i)(I_j - \pi_j)\right] \\ &= \sum_{i=1}^N \frac{x_i^2}{\pi_i^2} E(I_i - \pi_i)^2 + \sum_{i=1}^N \sum_{j \neq i}^N \frac{x_i}{\pi_i} \frac{x_j}{\pi_j} E(I_i I_j - I_i \pi_j - \pi_i I_j + \pi_i \pi_j) \\ &= \sum_{i=1}^N \frac{x_i^2}{\pi_i^2} D I_i + \sum_{i=1}^N \sum_{j \neq i}^N \frac{x_i}{\pi_i} \frac{x_j}{\pi_j} (E(I_i I_j) - \pi_j E I_i - \pi_i E I_j + \pi_i \pi_j) \\ &= \sum_{i=1}^N \frac{x_i^2}{\pi_i^2} \pi_i (1 - \pi_i) + \sum_{i=1}^N \sum_{j \neq i}^N \frac{x_i}{\pi_i} \frac{x_j}{\pi_j} (\pi_{ij} - \pi_j \pi_i - \pi_i \pi_j + \pi_i \pi_j) \\ &= \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i} \cdot x_i^2 + \sum_{i=1}^N \sum_{j \neq i}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} x_i x_j; \end{aligned}$$

$$\begin{aligned} 3) E(D(\hat{t}_{HT})) &= E\left[\sum_{i=1}^N \frac{1 - \pi_i}{\pi_i^2} \cdot x_i^2 + \sum_{i=1}^N \sum_{j \neq i}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{x_i x_j}{\pi_i \pi_j}\right] \\ &= E\left[\sum_{i=1}^N \frac{1 - \pi_i}{\pi_i^2} \cdot x_i^2 \cdot I_i + \sum_{i=1}^N \sum_{j \neq i}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{x_i x_j}{\pi_i \pi_j} I_i I_j\right] \\ &= \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i^2} \cdot x_i^2 E I_i + \sum_{i=1}^N \sum_{j \neq i}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{x_i x_j}{\pi_i \pi_j} E(I_i I_j) \\ &= \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i^2} \cdot x_i^2 \pi_i + \sum_{i=1}^N \sum_{j \neq i}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{x_i x_j}{\pi_i \pi_j} \cdot \pi_i \pi_j \\ &= D(\hat{t}_{HT}) \Rightarrow \text{јесте непристрасна.} \end{aligned}$$

Теорема 2: 1) Оцена $\bar{x}_{HT} = \frac{1}{N} \sum_{i=1}^N \frac{x_i}{\pi_i}$ је непристрасна оцена за \bar{x} ;

2) За дисперзију те оцене важи: $D(\bar{x}_{HT}) = \frac{1}{N^2} \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i} \cdot x_i^2 + \sum_{i=1}^N \sum_{j \neq i}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} x_i x_j$;

3) Непристрасна оцена те дисперзије је: $D(\hat{\bar{x}}_{HT}) = \frac{1}{N^2} \left[\sum_{i=1}^N \frac{1 - \pi_i}{\pi_i^2} \cdot x_i^2 + \sum_{i=1}^N \sum_{j \neq i}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{x_i x_j}{\pi_i \pi_j} \right]$;

Доказ: тривијално следи из Т1 и чињенице $\bar{x} = \frac{1}{N} \cdot t$.

Sen-Yates-Grundy - јеве оцене

Постоји боља оцена за $D(\hat{t}_{HT})$ од $\widehat{D(\hat{t}_{HT})}$. Показаћемо да је и она непристрасна.
Ова оцена се користи за фиксирано $\nu = n$ (тј. код узорка без понављања, и различитих)

Лема 1: 1) $\sum_{i=1}^N \pi_i = n$;
2) $\sum_{\substack{j=1 \\ j \neq i}}^N \pi_{ij} = (n-1) \pi_i$.

Доказ: 1) тривијално ($\sum \pi_i = \sum E I_i = E(\sum I_i) = E \nu = \nu = n$);

2) $\sum_{\substack{j=1 \\ j \neq i}}^N \pi_{ij} = \sum_{\substack{j=1 \\ j \neq i}}^N E(I_i I_j) = E(I_i \sum_{\substack{j=1 \\ j \neq i}}^N I_j) = E(I_i (n - I_i))$
 $= E(n I_i - I_i^2) = E(n I_i - I_i) = (n-1) E I_i = (n-1) \pi_i$.

Теорема 1: Непристрасна оцена за $D(\hat{t}_{HT})$ је: $D_{SYG}(\hat{t}_{HT}) = \sum_{i < j} \sum_{j \neq i} \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \left(\frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2$, $\pi_i, \pi_j, \pi_{ij} > 0$.

Доказ: Запишимо $D(\hat{t}_{HT})$ у алт. облику:

$$D(\hat{t}_{HT}) = \sum_{i=1}^N \frac{1 - \pi_i}{\pi_i} \cdot x_i^2 + \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} x_i x_j = \sum_{i=1}^N \frac{\pi_i (1 - \pi_i)}{\pi_i^2} \cdot x_i^2 - \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N \frac{\pi_i \pi_j - \pi_{ij}}{\pi_i \pi_j} x_i x_j$$

$$= \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{x_i}{\pi_i} \right)^2 - \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{x_i}{\pi_i} \cdot \frac{x_j}{\pi_j} \right)$$

$$= \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N (\pi_i \pi_j - \pi_{ij}) \left(\left(\frac{x_i}{\pi_i} \right)^2 + \left(\frac{x_j}{\pi_j} \right)^2 \right) - \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N (\pi_i \pi_j - \pi_{ij}) \cdot 2 \left(\frac{x_i}{\pi_i} \cdot \frac{x_j}{\pi_j} \right)$$

$$= \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2$$

$$(*) \sum_{\substack{i=1 \\ i \neq j}}^N (\pi_i \pi_j - \pi_{ij}) = \sum_{\substack{i=1 \\ i \neq j}}^N \pi_i \pi_j - \pi_i \sum_{\substack{j=1 \\ j \neq i}}^N \pi_j$$

$$\stackrel{(**)}{=} (n-1)\pi_i - \pi_i(n-\pi_i) = \pi_i(1-\pi_i)$$

$$E(D_{SYG}(\hat{t}_{HT})) = E \left[\sum_{i < j} \sum_{j \neq i} \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \left(\frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2 \right]$$

$$= \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \left(\frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2 E(I_i I_j)$$


$$= \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \left(\frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2 \pi_{ij} = \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N (\pi_i \pi_j - \pi_{ij}) \left(\frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2 \stackrel{\text{алт. облик}}{=} D(\hat{t}_{HT}) \Rightarrow \text{јесте непристрасна.}$$

$$(**) \text{ нпр. } n=4: \begin{array}{|c|c|c|c|} \hline 12 & 21 & 31 & 41 \\ \hline 13 & 23 & 32 & 42 \\ \hline 14 & 24 & 34 & 44 \\ \hline \end{array}$$

Остали су исти

Теорема 2: Непристрасна оцена за $D(\bar{x}_{HT})$ је: $D_{SYG}(\bar{x}_{HT}) = \frac{1}{N^2} \left[\sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N \frac{\pi_i \pi_j - \pi_{ij}}{\pi_{ij}} \left(\frac{x_i}{\pi_i} - \frac{x_j}{\pi_j} \right)^2 \right]$;

Доказ: тривијално из Т1

A hand-drawn teal diamond pattern, consisting of multiple overlapping diamond shapes, covers the entire page. The pattern is centered around the text.

II deo

11.

Стратификован узорак.

Веза укупне оцене и оцена по стратумима

* Имамо популацију коју делимо на делове/слојеве који се зову **стратуми**.

У оквиру једног стратума, вредности обележја треба да буду што хомогенија, а између стратума што хетерогенија.



$N = N_1 + \dots + N_L$: популација подељена на L стратума.

Проблем: Не знамо вредност обележја. Како направити стратуме онда?

Тражи се нека особина јединки коју знамо и да је у корелацији са вредностима обележја.

(нпр. обележје: плата, особина: ниво студија)

Бр. елем. у стратуму мора бити ≥ 2 .

Када направимо стратуме, вадио узорак обима $n = n_1 + \dots + n_L$ (n_h - независни)

Тако добијамо **стратификован узорак**.

Стратификован случајан узорак је онај стр. уз. код ког се узорци из сваког стратума ваде као ПСУ.

Напомена: Са страт. уз. се добијају боље оцене него код ПСУ, јер имамо одатне информације.

Постављају се питања: - на колико стратума поделити популацију, тј. које L узети?

- коју особину изабрати за поделу?

- како изабрати n_h -ове?

- * деф. N - обим популације;
- n - обим узорка
- L - број стратума;
- N_h - обим h -тог стратума;
- n_h - обим узорка из h -тог стратума.

деф. Обележја јединки h -тог стратума: X_{hi} , $i \in \{1, \dots, N_h\}$;

деф. Средња вредност обележја на h -том стратуму: $\bar{X}_h = \frac{\sum_i X_{hi}}{N_h}$;

деф. Укупна сума обележја на h -том стратуму: $t_h = \sum_i X_{hi}$;

деф. Пропорција јединки стратума са неким својством на h -том стратуму: $p_h = \frac{A_h}{N_h}$;

деф. Дисперзија на h -том стратуму: $S_h^2 = \frac{1}{N_h - 1} \sum_i (X_{hi} - \bar{X}_h)^2$.

деф. Узорачка средина на h -том стратуму: $\bar{X}_{n_h} = \frac{\sum_{i \in S_h} X_{hi}}{n_h}$;

деф. Оцена укупне суме на h -том стратуму: $\hat{t}_h = \sum_i X_{hi}$;

деф. Пропорција јединки стратума са својством p_h на узорку из h -тог стратума: $p_{n_h} = \frac{a_{n_h}}{n_h}$.

деф. Оцена дисперзије на h -том стратуму: $S_{n_h}^2 = \frac{1}{n_h - 1} \sum_i (X_{hi} - \bar{X}_{n_h})^2$.

Следеће две теореме су „опште“: важе за сваки стратификован узорак (не само случајни).

Теорема 1: Нека је у сваком стратуму \hat{t}_h непристрасна оцена за t_h .
и нека је у сваком стратуму $D(\hat{t}_h)$ непристрасна оцена за $D(t_h)$.

- 1) Оцена $\hat{t} = \sum_h \hat{t}_h$ је непристрасна оцена за t ;
- 2) За дисперзију те оцене важи: $D(\hat{t}) = \sum_h D(\hat{t}_h)$;
- 3) Непристрасна оцена те дисперзије је $D(\hat{t}) = \sum_h D(\hat{t}_h)$.

Доказ: 1) $E(\hat{t}) = E(\sum_h \hat{t}_h) = \sum_h E(\hat{t}_h) = \sum_h t_h = t \Rightarrow$ јесте непристрасна;

2) $D(\hat{t}) = D(\sum_h \hat{t}_h) \stackrel{\text{нез.}}{=} \sum_h D(\hat{t}_h)$;

3) $E(D(\hat{t})) = E(\sum_h D(\hat{t}_h)) = \sum_h E(D(\hat{t}_h)) = \sum_h D(t_h) = D(t) \Rightarrow$ јесте непристрасна.

Теорема 2: Нека је у сваком стратуму \hat{x}_h непристрасна оцена за \bar{x}_h
и нека је у сваком стратуму $D(\hat{x}_h)$ непристрасна оцена за $D(\bar{x}_h)$.

- 1) Оцена $\hat{x} = \frac{1}{N} \sum_h N_h \hat{x}_h$ је непристрасна оцена за \bar{x} ;
- 2) За дисперзију те оцене важи: $D(\hat{x}) = \frac{1}{N^2} \sum_h N_h^2 D(\hat{x}_h)$;
- 3) Непристрасна оцена те дисперзије је $D(\hat{x}) = \frac{1}{N^2} \sum_h N_h^2 D(\hat{x}_h)$;

Доказ: 1) $E(\hat{x}) = E(\frac{1}{N} \sum_h N_h \hat{x}_h) = \frac{1}{N} \sum_h N_h E(\hat{x}_h) = \frac{1}{N} \sum_h N_h \bar{x}_h = \frac{1}{N} \sum_h t_h = \frac{t}{N} = \bar{x} \Rightarrow$ јесте непристрасна;

2) $D(\hat{x}) = D(\frac{1}{N} \sum_h N_h \hat{x}_h) \stackrel{\text{нез.}}{=} \frac{1}{N^2} \sum_h N_h^2 D(\hat{x}_h)$;

3) $E(D(\hat{x})) = E(\frac{1}{N^2} \sum_h N_h^2 D(\hat{x}_h)) = \frac{1}{N^2} \sum_h N_h^2 E(D(\hat{x}_h)) = \frac{1}{N^2} \sum_h N_h^2 D(\bar{x}_h) = D(\bar{x}) \Rightarrow$ јесте непристрасна.

Стратификован случајан узорак са и без понављања

Сва тврђења у овом питању су последница / спец. случ. претходне две „опште“ теореме.

Следеће три теореме важе за узорак без понављања.

Теорема 1: Ако се из сваког стратума вади ПСУ без понављања:

- 1) Оцена $\hat{t}_{STR} = \sum_h^L N_h \cdot \bar{x}_{nh}$ је непристрасна оцена за t ;
- 2) За дисперзију те оцене важи: $D(\hat{t}_{STR}) = \sum_h^L \frac{N_h^2 \cdot s_{nh}^2}{n_h} \left(1 - \frac{n_h}{N_h}\right)$;
- 3) Непристрасна оцена те дисперзије је $\widehat{D}(\hat{t}_{STR}) = \sum_h^L \frac{N_h^2 \cdot s_{nh}^2}{n_h} \left(1 - \frac{n_h}{N_h}\right)$.

Теорема 2: Ако се из сваког стратума вади ПСУ без понављања:

- 1) Оцена $\bar{x}_{STR} = \frac{1}{N} \sum_h^L N_h \cdot \bar{x}_{nh}$ је непристрасна оцена за \bar{x} ;
- 2) За дисперзију те оцене важи: $D(\bar{x}_{STR}) = \frac{1}{N^2} \sum_h^L \frac{N_h^2 \cdot s_{nh}^2}{n_h} \left(1 - \frac{n_h}{N_h}\right)$;
- 3) Непристрасна оцена те дисперзије је $\widehat{D}(\bar{x}_{STR}) = \frac{1}{N^2} \sum_h^L \frac{N_h^2 \cdot s_{nh}^2}{n_h} \left(1 - \frac{n_h}{N_h}\right)$.

Теорема 3: Ако се из сваког стратума вади ПСУ без понављања:

- 1) Оцена $\hat{p}_{STR} = \frac{1}{N} \sum_h^L N_h \cdot p_{nh}$ је непристрасна оцена за p ;
- 2) За дисперзију те оцене важи: $D(\hat{p}_{STR}) = \frac{1}{N^2} \sum_h^L \frac{N_h^2 (N_h - n_h)}{n_h (N_h - 1)} p_{nh} (1 - p_{nh})$;
- 3) Непристрасна оцена те дисперзије је $\widehat{D}(\hat{p}_{STR}) = \frac{1}{N^2} \sum_h^L \frac{N_h^2 (N_h - n_h)}{n_h - 1} p_{nh} (1 - p_{nh})$.

Следеће три теореме важе за узорак са понављањем.

Теорема 4: Ако се из сваког стратума вади ПСУ са понављањем:

- 1) Оцена $\hat{t}_{STR} = \sum_h^L N_h \cdot \bar{x}_{nh}$ је непристрасна оцена за t ;
- 2) За дисперзију те оцене важи: $D(\hat{t}_{STR}) = \sum_h^L \frac{N_h(N_h-1)}{n_h} \cdot s_h^2$;
- 3) Непристрасна оцена те дисперзије је $\widehat{D}(\hat{t}_{STR}) = \sum_h^L N_h^2 \cdot \frac{s_{nh}^2}{n_h}$.

Теорема 5: Ако се из сваког стратума вади ПСУ са понављањем:

- 1) Оцена $\bar{x}_{STR} = \frac{1}{N} \sum_h^L N_h \cdot \bar{x}_{nh}$ је непристрасна оцена за \bar{x} ;
- 2) За дисперзију те оцене важи: $D(\bar{x}_{STR}) = \frac{1}{N^2} \sum_h^L \frac{N_h(N_h-1)}{n_h} \cdot s_h^2$;
- 3) Непристрасна оцена те дисперзије је $\widehat{D}(\bar{x}_{STR}) = \frac{1}{N^2} \sum_h^L N_h^2 \cdot \frac{s_{nh}^2}{n_h}$.

Теорема 6: Ако се из сваког стратума вади ПСУ са понављањем:

- 1) Оцена $\hat{p}_{STR} = \frac{1}{N} \sum_h^L N_h \cdot p_{nh}$ је непристрасна оцена за p ;
- 2) За дисперзију те оцене важи: $D(\hat{p}_{STR}) = \frac{1}{N^2} \sum_h^L N_h^2 \cdot \frac{p_h(1-p_h)}{n_h}$;
- 3) Непристрасна оцена те дисперзије је $\widehat{D}(\hat{p}_{STR}) = \frac{1}{N^2} \sum_h^L N_h^2 \cdot \frac{p_{nh}(1-p_{nh})}{n_h-1}$.

Напомена: У наставку подразумевамо ПСУ без понављања, осим ако не нагласимо другачије.

13.

Пропорционални избор обима узорка по стратумима

У наредна три питања гледамо избор обима по стратумима.

* За почетак, гледамо случај када су n_h -ови пропорционални N_h -овима: $\frac{n_1}{N_1} = \dots = \frac{n_L}{N_L} = \frac{n}{N} \Rightarrow n_h = \frac{n}{N} \cdot N_h$

Означимо са \hat{t}_{PROP} специјални случај за \hat{t}_{STR} под овим условима.

Од две оцене, боља је она са мањим MSE.

Како су и \hat{t}_{PROP} и \hat{t}_{STR} непристрасне, довољно је да им упоредимо дисперзије.

$$D(\hat{t}_{\text{PROP}}) \stackrel{\text{[17.2]}}{=} \sum_h N_h^2 \cdot \frac{s_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) \stackrel{n_h}{=} \sum_h \frac{N}{n} N_h \cdot s_h^2 \left(1 - \frac{n}{N}\right) = \frac{N}{n} \left(1 - \frac{n}{N}\right) \sum_h N_h s_h^2;$$

$$D(\hat{t}) \stackrel{\text{[17.2]}}{=} \frac{N^2 s^2}{n} \left(1 - \frac{n}{N}\right) = \frac{N}{n} \left(1 - \frac{n}{N}\right) \cdot N \cdot \frac{1}{N-1} \sum_i (x_i - \bar{x})^2$$

(*) Уместо да бројимо редом од x_1 до x_n , пребројавамо све у првом стратуму, па све у другом, па у трећем итд.

$$\stackrel{(*)}{=} \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{N}{N-1} \sum_h \sum_i^{N_h} (x_{hi} - \bar{x})^2$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{N}{N-1} \sum_h \sum_i^{N_h} (\underbrace{x_{hi} - \bar{x}_h}_{\text{green}} + \underbrace{\bar{x}_h - \bar{x}}_{\text{orange}})^2$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{N}{N-1} \sum_h \sum_i^{N_h} \left[(x_{hi} - \bar{x}_h)^2 + 2(x_{hi} - \bar{x}_h)(\bar{x}_h - \bar{x}) + (\bar{x}_h - \bar{x})^2 \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{N}{N-1} \left[\sum_h (N_h - 1) s_h^2 + 2 \sum_h (\bar{x}_h - \bar{x}) \sum_i^{N_h} (x_{hi} - \bar{x}_h) + \sum_h N_h \cdot (\bar{x}_h - \bar{x})^2 \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{N}{N-1} \left[\sum_h (N_h - 1) s_h^2 + 2 \sum_h (\bar{x}_h - \bar{x}) (N_h \bar{x}_h - N_h \bar{x}_h) + \sum_h N_h \cdot (\bar{x}_h - \bar{x})^2 \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{N}{N-1} \left[\sum_h (N_h - 1) s_h^2 + \sum_h N_h \cdot (\bar{x}_h - \bar{x})^2 \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \left(1 + \frac{1}{N-1}\right) \left[\sum_h N_h s_h^2 - \sum_h s_h^2 + \sum_h N_h (\bar{x}_h - \bar{x})^2 \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \left[\sum_h N_h s_h^2 - \sum_h s_h^2 + \frac{1}{N-1} \sum_h (N_h - 1) s_h^2 + \frac{N}{N-1} \sum_h N_h (\bar{x}_h - \bar{x})^2 \right]$$

$$= \frac{N}{n} \left(1 - \frac{n}{N}\right) \sum_h N_h s_h^2 + \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[-(N-1) \sum_h s_h^2 + \sum_h (N_h - 1) s_h^2 + N \sum_h N_h (\bar{x}_h - \bar{x})^2 \right]$$

$$= D(\hat{t}_{\text{PROP}}) + \frac{N}{n} \left(1 - \frac{n}{N}\right) \frac{1}{N-1} \left[N \sum_h N_h (\bar{x}_h - \bar{x})^2 - \sum_h (N - N_h) s_h^2 \right]$$

Пакле: $D(\hat{t}_{PROP}) \leq D(\hat{t})$ акко $N \sum_h^L N_h (\bar{X}_h - \bar{X})^2 \geq \sum_h^L (N - N_h) S_h^2 \stackrel{! : N}{\Rightarrow} \sum_h^L N_h (\bar{X}_h - \bar{X})^2 \geq \sum_h^L \overbrace{\left(1 - \frac{N_h}{N}\right)}^{\text{фикс.}} S_h^2$

↑ што веће
што веће

што мање

Закључак:

\hat{t}_{PROP} ће бити боља од \hat{t} када је:

→ дисперзија у сваком стратуму **што мање**

(што хомогенији у оквиру стратума)

→ ср. вр. на стратумима **што различитије** од ср. вр. целе популације

(што хетерогеније између стратума)

→ N_h да буде **што веће**.

14.

Оптимални избор обима узорка по стратумима

Када су дисперзије стратума s_h^2 приближно једнаке, пропорц. избор n_h -ова ради добар посао. Ако ипак s_h^2 -ови варирају - користимо тзв. оптимални избор - тиме се сада бавимо.

* Видели смо шта се дешава ако су величине узорака из сваког стратума пропорционалне. Сада нас занима: колики узорак узети из сваког од стратума како би оцена била што боља?

За непристрасне оцене важи да су боље ако имају мању дисперзију.

Имамо проблем условног екстремума: за унапред задате трошкове, тражимо најмању дисперзију. (самим тим и најбољу оцену)

$$(*) \text{ Подсетимо се: } D(\hat{t}_{STR}) \stackrel{[2] \text{ Th. 2}}{=} \sum_h \frac{N_h^2}{n_h} \cdot s_h^2 \left(1 - \frac{n_h}{N_h}\right) = \sum_h \frac{N_h^2}{n_h} \cdot s_h^2 - \sum_h N_h \cdot s_h^2$$

деф. (**)
 $C = C_0 + \sum_{h=1}^L C_h \cdot n_h$, где: C - укупни трошкови;
 C_0 - фиксни трошкови узорковања;

C_h - трошкови по јединки h -тог стратума, променљиво - познато;
 n_h - обим узорка из h -тог стратум - непознато.

Циљ: тражимо минимум дисперзије, тј. функције (*) уз услов о трошковима (**).

То радимо тако што правимо Лагранжову ф-ју за тражење условног екстремума за ф-ју више променљивих. При томе, тражимо непознате n_1, \dots, n_L .

$$L(n_1, \dots, n_L, \lambda) = \sum_h^L \frac{N_h^2}{n_h} \cdot \Delta_h^2 - \sum_h^L N_h \cdot \Delta_h^2 + \lambda \left(\sum_h^L C_h \cdot n_h + c_0 - C \right)$$

$$= \sum_h^L \frac{N_h^2}{n_h} \cdot \Delta_h^2 \left(1 - \frac{n_h}{N_h} \right) + \lambda \left(\sum_h^L C_h \cdot n_h + c_0 - C \right)$$

(Φ -ја чији се минимум тражи)
+ λ · услов

$$\left. \begin{aligned} \frac{\partial L}{\partial n_h} &= -\frac{N_h^2 \Delta_h^2}{n_h^2} + \lambda C_h = 0, \quad h=1, \dots, L \\ \frac{\partial L}{\partial \lambda} &= \sum_h^L C_h n_h + c_0 - C = 0 \end{aligned} \right\} \text{систем } (L+1) \text{ једначина}$$

$$\lambda C_h = \frac{N_h^2 \Delta_h^2}{n_h^2} \Rightarrow n_h = \frac{N_h \Delta_h}{\sqrt{\lambda} \sqrt{C_h}}$$

$$\sum_h^L C_h n_h = \sum_h^L C_h \frac{N_h \Delta_h}{\sqrt{\lambda} \sqrt{C_h}} = C - c_0 \Rightarrow \frac{1}{\sqrt{\lambda}} \cdot \sum_h^L \sqrt{C_h} N_h \Delta_h = C - c_0 \Rightarrow \sqrt{\lambda} = \frac{\sum_h^L \sqrt{C_h} N_h \Delta_h}{C - c_0}$$

бројан ће бити k
пошто је k заједно

$$\Rightarrow n_h = \frac{N_h \Delta_h}{\frac{\sum_k^L \sqrt{C_k} N_k \Delta_k}{C - c_0} \sqrt{C_h}} \Rightarrow n_h = \frac{(C - c_0) \frac{N_h \Delta_h}{\sqrt{C_h}}}{\sum_k^L \sqrt{C_k} N_k \Delta_k} \Rightarrow n = \frac{(C - c_0) \sum_k^L \frac{N_k \Delta_k}{\sqrt{C_k}}}{\sum_k^L \sqrt{C_k} N_k \Delta_k}$$

Ова стационарна тачка је кандидат за минимум. То морамо да проверимо.
Зато рачунамо и друге парцијалне изводе (који иду у матрицу Φ):

$$\left. \begin{aligned} \frac{\partial^2 L}{\partial n_h^2} &= \frac{2 N_h^2 \Delta_h^2}{n_h^3} \\ \frac{\partial^2 L}{\partial n_i \partial n_j} &= 0, \quad i \neq j \end{aligned} \right\} \Rightarrow \Phi = \begin{bmatrix} \frac{2 N_1^2 \Delta_1^2}{n_1^3} & & 0 \\ & \ddots & \\ 0 & & \frac{2 N_L^2 \Delta_L^2}{n_L^3} \end{bmatrix} \quad \left(\begin{array}{l} \Phi \text{ поз. дефинитна} \Rightarrow \text{минимум;} \\ \Phi \text{ нег. дефинитна} \Rightarrow \text{максимум.} \end{array} \right)$$

Како је Φ позитивно дефинитна \Rightarrow у тој тачки јесте минимум.

Напомена: Може се десити $n_h > N_h$.

Тада се за n_h проглашава баш $n_h = N_h$, а онда поновимо процес, али без тог стратума.

15. Неутан-ов избор обима узорка по стратумима

* Посматрамо спец. случај оптималног узорковања, када је $c_1 = \dots = c_L = c$.
Такав избор зове се **Нејманов избор узорка**.

Изводимо цео поступак поново.

$$C = c_0 + \sum_{h=1}^L c_h n_h, \quad \text{па кад наметнемо услов: } C = c_0 + c \cdot n.$$

Како су трошкови унапред познати, то значи и да ће n бити познато (тј. фиксно).

Опет тражимо минимум по n_h за $D(\hat{t}_{STR}) = \sum_{h=1}^L \frac{N_h^2}{n_h} s_h^2 \left(1 - \frac{n_h}{N_h}\right)$, уз услов $n_1 + \dots + n_L = n$.

$$L(n_1, \dots, n_L, \lambda) = \sum_{h=1}^L \frac{N_h^2}{n_h} s_h^2 \left(1 - \frac{n_h}{N_h}\right) + \lambda(n_1 + \dots + n_L - n) = \sum_{h=1}^L \frac{N_h^2}{n_h} s_h^2 - \sum_{h=1}^L N_h \cdot s_h^2 + \lambda(n_1 + \dots + n_L - n)$$

$$\frac{\partial L}{\partial n_h} = -\frac{N_h^2 s_h^2}{n_h^2} + \lambda = 0, \quad h=1, \dots, L \quad \Rightarrow \quad n_h = \frac{N_h s_h}{\sqrt{\lambda}}$$

$$\frac{\partial L}{\partial \lambda} = n_1 + \dots + n_L - n = 0 \quad \Rightarrow \quad n = \sum_{h=1}^L \frac{N_h s_h}{\sqrt{\lambda}} \quad \Rightarrow \quad \sqrt{\lambda} = \frac{\sum_{k=1}^L N_k s_k}{n} \quad \Rightarrow \quad n_h = \frac{N_h s_h}{\sum_{k=1}^L N_k s_k} \cdot n$$

Нашли смо стационарну тачку (n_h), па проверавамо да ли је минимум:

$$\left. \begin{aligned} \frac{\partial^2 L}{\partial n_h^2} &= \frac{2 N_h^2 s_h^2}{n_h^3} \\ \frac{\partial^2 L}{\partial n_i \partial n_j} &= 0, \quad i \neq j \end{aligned} \right\} \Rightarrow \quad \phi = \begin{bmatrix} \frac{2 N_1^2 s_1^2}{n_1^3} & & 0 \\ & \dots & \\ 0 & & \frac{2 N_L^2 s_L^2}{n_L^3} \end{bmatrix}$$

Како је ϕ позитивно дефинитна \Rightarrow у тој тачки јесте минимум.

$$\text{Када то уврстимо: } D(\hat{t}_{STR}) = \sum_{h=1}^L \frac{N_h^2}{n_h} s_h^2 - \sum_{h=1}^L N_h \cdot s_h^2 = \sum_{h=1}^L \frac{N_h^2 s_h^2}{\frac{N_h s_h}{\sum_{k=1}^L N_k s_k} n} - \sum_{h=1}^L N_h \cdot s_h^2$$

$$D(\hat{t}_{NEU}) = \frac{1}{n} \left(\sum_{h=1}^L N_h s_h \right)^2 - \sum_{h=1}^L N_h s_h^2.$$

* **Проблем:** Да бисмо знали колики узорак бирамо, морамо да знамо s_h -ове.

Решење: Исто као онај пут: нпр. прелиминарно истраживање.

* Како се све ове оцене пореде међусобно?


\hat{t} , \hat{t}_{PROP} , \hat{t}_{NEU} могу да се пореде, зато што је n унапред познато.

Оцена добијена оптималним узорковањем не може да се пореди са њима, јер n није познато.

Лема 1: $D(\hat{t}_{NEU}) \leq D(\hat{t}_{PROP})$.

$$\begin{aligned}
 \text{Доказ: } D(\hat{t}_{PROP}) - D(\hat{t}_{NEU}) &= \frac{N-n}{n} \sum_h N_h \Delta_h^2 - \frac{1}{n} \left(\sum_h N_h \Delta_h \right)^2 + \sum_h N_h \Delta_h^2 \\
 &= \frac{N}{n} \sum_h N_h \Delta_h^2 - \frac{N}{n} \cdot \frac{1}{N} \left(\sum_h N_h \Delta_h \right)^2 \\
 &= \frac{N}{n} \left[\sum_h N_h \Delta_h^2 - \frac{2}{N} \left(\sum_h N_h \Delta_h \right)^2 + \frac{1}{N} \left(\sum_h N_h \Delta_h \right)^2 \right] \\
 &= \frac{N}{n} \left[\sum_h N_h \Delta_h^2 - \frac{2}{N} \left(\sum_h N_h \Delta_h \right) \left(\sum_k N_k \Delta_k \right) + \frac{\sum_h N_h}{N^2} \left(\sum_k N_k \Delta_k \right)^2 \right] \\
 &= \frac{N}{n} \sum_h N_h \left[\Delta_h^2 - 2\Delta_h \cdot \frac{1}{N} \sum_k N_k \Delta_k + \left(\frac{1}{N} \sum_k N_k \Delta_k \right)^2 \right] \\
 &= \frac{N}{n} \sum_h N_h \left[\Delta_h - \frac{1}{N} \sum_k N_k \Delta_k \right]^2 \geq 0
 \end{aligned}$$

Закључак: $D(\hat{t}_{NEU}) \leq D(\hat{t}_{PROP}) \leq D(\hat{t})$
 \uparrow увек \uparrow углавном (хомог. + хетерог.)

A decorative border made of hand-drawn teal lines forming a diamond or lattice pattern, surrounding the central text.

III

neo

Количничко оцењивање

* За почетак, објаснимо шта је количничко оцењивање.

Већ код стратификовања смо имали неку помоћну променљиву / обележје.

Боље оцене можемо добити на два начина:

а) узоровање вршимо на основу те помоћне променљиве - лакше, стратификован узорак;

б) на основу те помоћне променљиве правимо боље оцене у оквиру ПСУ.

Ово се зове **количничко оцењивање** и **линеарно - регресионо оцењивање**.

Код количничког оцењивања, имамо главно обележје x и помоћно обележје y .

Оцена на основу количничког оцењивања ће бити боља ако:

→ x, y су у корелацији;

→ ако из $y=0$ следи и $x=0$.

Пример: x - бр. животиња у шуми
 y - површина шуме

Контра пример: x - принос житарице
 y - загађеност ваздуха

Јесу у корелацији, али ако је $y=0$, принос расте.

21.

Количничко оцењивање на основу ПСУ без понављања

Имамо узорак: $(x_1, y_1), \dots, (x_n, y_n)$ обима n .

Напомена: За количничко оцењивање, увек морамо да знамо колико је $t_y = \sum_{i=1}^n y_i$

деф. Количник популације је $R := \frac{\sum_{i=1}^N x_i}{\sum_{i=1}^N y_i} = \frac{t}{t_y} \stackrel{t=N\bar{x}}{=} \frac{\bar{x}}{\bar{y}}$

деф. Количничка оцена за R је **узорачки количник** $\hat{R} := \frac{\sum_{i=1}^n x_i}{\sum_{i=1}^n y_i} = \frac{\bar{x}_n}{\bar{y}_n}$

деф. Количничка оцена за t је $\hat{t}_R = \hat{R} \cdot t_y$. (t_y је познато)

Количничка оцена за \bar{x} је $\bar{X}_R = \hat{R} \cdot \bar{y}$. ($\bar{y} = \frac{t_y}{N}$ је познато)

Када се користи количничко оцењивање?

- за оцењивање количника R ;
- за оцењивање t када N није познато, али t_y јесте познато;
- некад добијемо боље оцене за t и \bar{x} . [22]

* Испитујемо пристрасност ових оцена. Испоставиће се да су асимптотски непристрасне. Због тога, за велико n , за меру њиховог квалитета се опет користи дисперзија, уместо MSE.

Теорема 1: \hat{R} је асимптотски непристрасна.

Доказ: $\hat{R} - R = \frac{\bar{x}_n}{\bar{y}_n} - R = \frac{\bar{x}_n - \bar{y}_n \cdot R}{\bar{y}_n} \approx \frac{\bar{x}_n - R \cdot \bar{y}_n}{\bar{y}}$ за велико n : $\bar{y}_n \approx \bar{y}$

$$E(\hat{R} - R) \approx E\left(\frac{\bar{x}_n - R \cdot \bar{y}_n}{\bar{y}}\right) = \frac{1}{\bar{y}} E(\bar{x}_n - R \cdot \bar{y}_n) = \frac{1}{\bar{y}} [E(\bar{x}_n) - R \cdot E(\bar{y}_n)] \stackrel{псу}{=} \frac{\bar{x} - R \bar{y}}{\bar{y}} = 0$$

$$\Rightarrow E(\hat{R}) \approx R.$$

Последица: \hat{t}_R и \bar{X}_R су асимптотски непристрасне оцене. (јер $t_y = \text{const}$)

Теорема 2: За велико n , важе следеће апроксимације:

$$\begin{aligned}
 & \text{1) } \text{MSE}(\hat{R}) \approx D(\hat{R}) = \frac{(1-\frac{n}{N})}{n\bar{y}^2} \cdot \frac{1}{N-1} \cdot \sum_{i=1}^N (x_i - Ry_i)^2; \\
 & \text{2) } \text{MSE}(\hat{t}_R) \approx D(\hat{t}_R) = \frac{N^2(1-\frac{n}{N})}{n} \cdot \frac{1}{N-1} \cdot \sum_{i=1}^N (x_i - Ry_i)^2; \quad (\text{алтернативни облик у [22]}) \\
 & \text{3) } \text{MSE}(\bar{X}_R) \approx D(\bar{X}_R) = \frac{(1-\frac{n}{N})}{n} \cdot \frac{1}{N-1} \cdot \sum_{i=1}^N (x_i - Ry_i)^2.
 \end{aligned}$$

Доказ: 1) Због Т1, јасно је да важи $\text{MSE} \approx D$.

$$\begin{aligned}
 \text{MSE}(\hat{R}) &= E(\hat{R} - R)^2 \approx E\left(\frac{\bar{X}_n - R \cdot \bar{Y}_n}{\bar{y}}\right)^2 = \frac{1}{\bar{y}^2} [E(\bar{X}_n - R \cdot \bar{Y}_n)^2] \\
 &= \frac{1}{\bar{y}^2} [E(\bar{X}_n - R \cdot \bar{Y}_n)^2 - \underbrace{[E(\bar{X}_n - R \cdot \bar{Y}_n)]^2}_{=0, \text{ па може}}] = \frac{1}{\bar{y}^2} D(\bar{X}_n - R \cdot \bar{Y}_n) = \frac{1}{\bar{y}^2} D(\bar{d}_n) \\
 &\stackrel{\text{Т1}}{=} \frac{1}{\bar{y}^2} (1 - \frac{n}{N}) \frac{\hat{\Delta}_J^2}{n} = \frac{1}{\bar{y}^2} (1 - \frac{n}{N}) \cdot \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N (d_i - \bar{d})^2 = \frac{(1-\frac{n}{N})}{n\bar{y}^2} \cdot \frac{1}{N-1} \cdot \sum_{i=1}^N (x_i - Ry_i)^2;
 \end{aligned}$$

ОЗНАКА

$$\begin{aligned}
 \bar{X}_n - R \cdot \bar{Y}_n &= \frac{1}{n} \sum_{i=1}^n x_i - R \cdot \frac{1}{n} \sum_{i=1}^n y_i \\
 &= \frac{1}{n} \sum_{i=1}^n (x_i - R \cdot y_i) = \bar{d}_n
 \end{aligned}$$

2) тривијално (стрелица);

3) тривијално (стрелица).

Напомена: За ове дисперзије предлажемо следеће оцене:

$$\widehat{D}(\hat{t}_R) := \frac{N^2(1-\frac{n}{N})}{n} \cdot \frac{1}{n-1} \cdot \sum_{i \in S} (x_i - \hat{R}y_i)^2;$$

$$\widehat{D}(\bar{X}_R) := \frac{(1-\frac{n}{N})}{n} \cdot \frac{1}{n-1} \cdot \sum_{i \in S} (x_i - \hat{R}y_i)^2;$$

Ове оцене су асимптотски непристрасне. (зато што код ПСУ без пон. важи $E(\Delta_n^2) = \Delta^2$, па ће важити $E(\Delta_{n,J}^2) = \Delta_J^2$.)

22.

Упоредивање количничке оцене и оцене на основу ПСУ без понављања

* Подсетимо се: Коefицијент корелације мери лин. зав. два обележја ($|r| \approx 1 \Rightarrow y \approx kx + n$)

$$r_{x,y} := \frac{E[(X-EX)(Y-EY)]}{\sqrt{DX} \cdot \sqrt{DY}} = \frac{\frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{N-1}{N} \cdot s_x^2} \cdot \sqrt{\frac{N-1}{N} \cdot s_y^2}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{(N-1) s_x s_y}$$

$$\begin{aligned} x: \left(\frac{x_1}{N}, \dots, \frac{x_N}{N} \right) &\Rightarrow EX = \frac{1}{N} \sum x_i, & EX^2 &= \frac{1}{N} \sum x_i^2 \\ \hookrightarrow \text{псу б.п.} && DX &= EX^2 - (EX)^2 = \frac{1}{N} \sum x_i^2 - \bar{x}^2 = \frac{1}{N} (\sum x_i^2 - N\bar{x}^2) = \frac{N-1}{N} s_x^2 \end{aligned}$$

(екв) из [1]

* Поредимо оцене \hat{t}_R и \hat{t} (псу без понављања)

Како су обе непристрасне (\hat{t}_R асимптотски), поредимо им дисперзије:

$$\rightarrow D(\hat{t}) \stackrel{[1]}{=} N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} s_x^2$$

$$\begin{aligned} \rightarrow D(\hat{t}_R) &\stackrel{[2]}{=} \frac{N^2 \left(1 - \frac{n}{N}\right)}{n} \cdot \frac{1}{N-1} \cdot \sum (x_i - R y_i)^2 \stackrel{\text{ПРАВИЛО АНТ. ОБРАЗ}}{=} \frac{N^2 \left(1 - \frac{n}{N}\right)}{n} \cdot \frac{1}{N-1} \cdot \sum (x_i - R y_i - (\bar{x} - R \bar{y}))^2 \\ &= \frac{N^2 \left(1 - \frac{n}{N}\right)}{n} \cdot \frac{1}{N-1} \cdot \sum (x_i - \bar{x} - R(y_i - \bar{y}))^2 \\ &= \frac{N^2 \left(1 - \frac{n}{N}\right)}{n} \cdot \frac{1}{N-1} \cdot \left(\sum (x_i - \bar{x})^2 + \sum R^2 (y_i - \bar{y})^2 - 2R \sum (x_i - \bar{x})(y_i - \bar{y}) \right) \\ &= \frac{N^2 \left(1 - \frac{n}{N}\right)}{n} \left(\underline{s_x^2} + R^2 \underline{s_y^2} - 2R \rho \underline{s_x s_y} \right) \end{aligned}$$

погледати деф. ρ

$$\sum (x_i - \bar{x})(y_i - \bar{y}) = \rho \cdot (N-1) \cdot s_x s_y$$

$$\hat{t}_R \text{ боља од } \hat{t} \Leftrightarrow D(\hat{t}_R) - D(\hat{t}) = N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} (R^2 s_y^2 - 2R \rho s_x s_y) < 0$$

$$\Leftrightarrow R^2 s_y^2 - 2R \rho s_x s_y < 0$$

$$\stackrel{R > 0}{\Leftrightarrow} \rho > \frac{1}{2} R \frac{s_y}{s_x}$$

$$\Leftrightarrow \rho > \frac{1}{2} \cdot \frac{CV_y}{CV_x}, \quad CV_x := \frac{s_x}{\bar{x}} \text{ је коefицијент варијације.}$$

Закључак: \hat{t}_R је боља од \hat{t} када је ρ што веће, тј. што је већа лин. зав. између x и y .

23.

Количничко оцењивање на основу узорка са неједнаким вероватноћама избора

Користимо Horvitz - Thompson - ове оцене.^[2]

деф. За оцену R , користимо $\hat{R}_{HT} := \frac{\hat{t}_{HT}(x)}{\hat{t}_{HT}(y)} = \frac{\sum \frac{x_i}{\pi_i}}{\sum \frac{y_i}{\pi_i}}$. (писаћемо опет само \hat{R})

деф. $\hat{t}_R := \hat{R}_{HT} \cdot t_y$.

$\bar{x}_R := \hat{R}_{HT} \cdot \bar{y}$.

Радио аналогно као у [21]:

Теорема 1: \hat{R}_{HT} је асимптотски непристрасна.

Доказ: $\hat{R} - R = \frac{\hat{t}_{HT}(x)}{\hat{t}_{HT}(y)} - R = \frac{\hat{t}_{HT}(x) - R \cdot \hat{t}_{HT}(y)}{\hat{t}_{HT}(y)} \approx \frac{\hat{t}_{HT}(x) - R \cdot \hat{t}_{HT}(y)}{t_y}$

за велико n : $\hat{t}_{HT}(y) \approx t_y$

$$E(\hat{R} - R) \approx E\left(\frac{\hat{t}_{HT}(x) - R \cdot \hat{t}_{HT}(y)}{t_y}\right) = \frac{1}{t_y} [E(\hat{t}_{HT}(x)) - R \cdot E(\hat{t}_{HT}(y))] \stackrel{[2][4]}{=} \frac{1}{t_y} (t_x - R t_y) = 0.$$

Последица: \hat{t}_R и \bar{x}_R су асимптотски непристрасне оцене.

Теорема 2: За велико n , важе следеће апроксимације:

$$\begin{aligned}
 & \cdot \frac{1}{t_y^2} \left\{ \begin{aligned}
 1) \text{ MSE}(\hat{R}_{HT}) &\approx D(\hat{R}_{HT}) = \frac{1}{t_y^2} \left[\sum_{i=4}^N \frac{1-\pi_i}{\pi_i} \cdot (x_i - R \cdot y_i)^2 + \sum_{i=4}^N \sum_{j=i+1}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} (x_i - R \cdot y_i)(x_j - R \cdot y_j) \right]; \\
 2) \text{ MSE}(\hat{t}_R) &\approx D(\hat{t}_R) = \sum_{i=4}^N \frac{1-\pi_i}{\pi_i} \cdot (x_i - R \cdot y_i)^2 + \sum_{i=4}^N \sum_{j=i+1}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} (x_i - R \cdot y_i)(x_j - R \cdot y_j); \\
 3) \text{ MSE}(\bar{X}_R) &\approx D(\bar{X}_R) = \frac{1}{N^2} \left[\sum_{i=4}^N \frac{1-\pi_i}{\pi_i} \cdot (x_i - R \cdot y_i)^2 + \sum_{i=4}^N \sum_{j=i+1}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} (x_i - R \cdot y_i)(x_j - R \cdot y_j) \right].
 \end{aligned} \right.
 \end{aligned}$$

Доказ: 1) Због Т1, јасно је да важи $\text{MSE} \approx D$.

$$\begin{aligned}
 \text{MSE}(\hat{R}) &= E(\hat{R} - R)^2 \approx E\left(\frac{\hat{t}_{HT}(x) - R \cdot \hat{t}_{HT}(y)}{t_y}\right)^2 = \frac{1}{t_y^2} \left[E(\hat{t}_{HT}(x) - R \hat{t}_{HT}(y))^2 \right] \\
 &= \frac{1}{t_y^2} \left[E(\hat{t}_{HT}(x) - R \hat{t}_{HT}(y))^2 - \underbrace{(E(\hat{t}_{HT}(x) - R \hat{t}_{HT}(y)))^2}_{=0, \text{ по моме}} \right] = \frac{1}{t_y^2} D(\hat{t}_{HT}(x) - R \hat{t}_{HT}(y)) \\
 &= \frac{1}{t_y^2} D\left(\sum_{i=4}^N \frac{x_i}{\pi_i} - R \cdot \sum_{i=4}^N \frac{y_i}{\pi_i}\right) = \frac{1}{t_y^2} \cdot D\left(\sum_{i=4}^N \frac{x_i - R y_i}{\pi_i}\right) = \frac{1}{t_y^2} D(\hat{t}_{HT}(x - R \cdot y)) \\
 &\stackrel{\text{Т1}}{=} \frac{1}{t_y^2} \left[\sum_{i=4}^N \frac{1-\pi_i}{\pi_i} \cdot (x_i - R \cdot y_i)^2 + \sum_{i=4}^N \sum_{j=i+1}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} (x_i - R \cdot y_i)(x_j - R \cdot y_j) \right]
 \end{aligned}$$

2) тривијално (стрелица);

3) тривијално (стрелица).

Напомена: За ове дисперзије предлажемо следеће оцене:

$$\begin{aligned}
 & \cdot \frac{1}{N^2} \left\{ \begin{aligned}
 \hat{D}(\hat{t}_R) &:= \sum_{i=4}^N \frac{1-\pi_i}{\pi_i^2} \cdot (x_i - \hat{R} \cdot y_i)^2 + \sum_{i=4}^N \sum_{j=i+1}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{(x_i - \hat{R} \cdot y_i)(x_j - \hat{R} \cdot y_j)}{\pi_{ij}}; \\
 \hat{D}(\bar{X}_R) &:= \frac{1}{N^2} \left[\sum_{i=4}^N \frac{1-\pi_i}{\pi_i^2} \cdot (x_i - \hat{R} \cdot y_i)^2 + \sum_{i=4}^N \sum_{j=i+1}^N \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{(x_i - \hat{R} \cdot y_i)(x_j - \hat{R} \cdot y_j)}{\pi_{ij}} \right].
 \end{aligned} \right.
 \end{aligned}$$

24.

Количничко оцењивање на основу стратификованог узорка без понављања комбинована оцена

деф. За оцену R , користимо $\hat{R}_c := \frac{\hat{t}_{STR}(x)}{\hat{t}_{STR}(y)} = \frac{\sum_{L} N_h \cdot x_{nh}}{\sum_{L} N_h \cdot y_{nh}} = \frac{\bar{x}_{STR}}{\bar{y}_{STR}}$. (писаћемо опет само \hat{R})

деф. $\hat{t}_{RC} := \hat{R}_c \cdot t_y$.

$\bar{x}_{RC} := \hat{R}_c \cdot \bar{y}$.

Опет исто.

Теорема 1: \hat{R}_c је асимптотски непристрасна.

за велико n : $\bar{y}_{STR} \approx \bar{y}$

Доказ: $\hat{R} - R = \frac{\bar{x}_{STR}}{\bar{y}_{STR}} - R = \frac{\bar{x}_{STR} - R \cdot \bar{y}_{STR}}{\bar{y}_{STR}} \approx \frac{\bar{x}_{STR} - R \cdot \bar{y}}{\bar{y}}$

$$E(\hat{R} - R) \approx E\left(\frac{\bar{x}_{STR} - R \cdot \bar{y}_{STR}}{\bar{y}}\right) = \frac{1}{\bar{y}} [E(\bar{x}_{STR}) - R \cdot E(\bar{y}_{STR})] \stackrel{E2T2}{=} \frac{1}{\bar{y}} (\bar{x} - R \cdot \bar{y}) = 0.$$

Последица: \hat{t}_{RC} и \bar{x}_{RC} су асимптотски непристрасне оцене.

Теорема 2: За велико n , важе следеће апроксимације:

$$\begin{aligned}
 & \cdot \frac{1}{t_y^2} \left. \begin{aligned}
 1) \text{ MSE}(\hat{R}_c) &\approx D(\hat{R}_c) = \frac{1}{t_y^2} \left[\sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_h^2(x) + R^2 \Delta_h^2(y) - 2R \rho_h \Delta_h(x) \Delta_h(y)) \right]; \\
 2) \text{ MSE}(\hat{t}_{RC}) &\approx D(\hat{t}_{RC}) = \sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_h^2(x) + R^2 \Delta_h^2(y) - 2R \rho_h \Delta_h(x) \Delta_h(y)); \\
 3) \text{ MSE}(\bar{X}_{RC}) &\approx D(\bar{X}_{RC}) = \frac{1}{N^2} \left[\sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_h^2(x) + R^2 \Delta_h^2(y) - 2R \rho_h \Delta_h(x) \Delta_h(y)) \right].
 \end{aligned} \right. \\
 & \cdot \frac{1}{N^2}
 \end{aligned}$$

Доказ: 1) Због Т1, јасно је да важи $\text{MSE} \approx D$.

$$\begin{aligned}
 \text{MSE}(\hat{R}) &= E(\hat{R} - R)^2 \approx E\left(\frac{\bar{X}_{STR} - R \cdot \bar{Y}_{STR}}{\bar{y}}\right)^2 = \frac{1}{\bar{y}^2} \left[E(\bar{X}_{STR} - R \cdot \bar{Y}_{STR})^2 \right] \\
 &= \frac{1}{\bar{y}^2} \left[E(\bar{X}_{STR} - R \cdot \bar{Y}_{STR})^2 - \underbrace{(E(\bar{X}_{STR} - R \cdot \bar{Y}_{STR}))^2}_{=0, \text{ на моме}} \right] = \frac{1}{\bar{y}^2} D(\bar{X}_{STR} - R \cdot \bar{Y}_{STR}) \\
 &= \frac{1}{\bar{y}^2} D\left(\frac{1}{N} \sum_h N_h \bar{X}_{nh} - R \cdot \frac{1}{N} \sum_h N_h \bar{Y}_{nh}\right) = \frac{1}{\bar{y}^2} D\left(\frac{1}{N} \sum_h N_h \frac{\sum_{i \in S} (X_{hi} - R Y_{hi})}{n_h}\right) = \frac{1}{\bar{y}^2} D((X-RY)_{STR}) \quad \nearrow \text{ново оделење} \\
 &\stackrel{12)}{=} \frac{1}{\bar{y}^2} \cdot \frac{1}{N^2} \sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{1}{N_h - 1} \sum_i^{N_h} (X_{hi} - R Y_{hi} - (\bar{X}_h - R \bar{Y}_h))^2 \\
 &= \frac{1}{t_y^2} \sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) \frac{1}{N_h - 1} \sum_i^{N_h} ((X_{hi} - \bar{X}_h) - R(Y_{hi} - \bar{Y}_h))^2 \\
 &= \frac{1}{t_y^2} \sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) \frac{1}{N_h - 1} \sum_i^{N_h} ((X_{hi} - \bar{X}_h)^2 + R^2 (Y_{hi} - \bar{Y}_h)^2 - 2R(X_{hi} - \bar{X}_h)(Y_{hi} - \bar{Y}_h)) \\
 &= \frac{1}{t_y^2} \left[\sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_h^2(x) + R^2 \Delta_h^2(y) - 2R \rho_h \Delta_h(x) \Delta_h(y)) \right]
 \end{aligned}$$

2) тривијално (стрелица);

3) тривијално (стрелица).

Напомена: Предлажемо следећу оцену:

$$\begin{aligned}
 & \cdot \frac{1}{N^2} \left. \begin{aligned}
 \hat{D}(\hat{t}_{RC}) &:= \sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_{nh}^2(x) + \hat{R}^2 \Delta_{nh}^2(y) - 2\hat{R} \rho_h \Delta_{nh}(x) \Delta_{nh}(y)); \\
 \hat{D}(\bar{X}_{RC}) &:= \frac{1}{N^2} \left[\sum_h \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_{nh}^2(x) + \hat{R}^2 \Delta_{nh}^2(y) - 2\hat{R} \rho_h \Delta_{nh}(x) \Delta_{nh}(y)) \right].
 \end{aligned} \right.
 \end{aligned}$$

25.

Количничко оцењивање на основу стратификованог узорка без понављања посебна по стратумима оцена

* На сваком стратуму оцењујемо укупну суму и онда их саберемо. (зато иде \sum_h^L)

деф. $\hat{R}_h := \frac{\hat{t}_h(x)}{\hat{t}_h(y)} = \frac{\bar{x}_{n_h}}{\bar{y}_{n_h}}$; - по [24], ово је асимптотски непристрасно (стратум = псу)

деф. $\hat{t}_{RS} := \sum_h^L \hat{R}_h \cdot t_h(y) = \sum_h^L \frac{\sum_{i \in S_h} x_{hi}}{\sum_{i \in S_h} y_{hi}} \cdot t_h(y)$

$\bar{x}_{RS} := \frac{1}{N} \cdot \hat{t}_{RS}$

Теорема 1: \hat{t}_{RS} је асимптотски непристрасна.

Доказ: $E(\hat{t}_{RS}) = \sum_h^L t_h(y) \cdot E(\hat{R}_h) \approx \sum_h^L t_h(y) \cdot R_h = \sum_h^L t_h(x) = t$

Теорема 2: За велико n , важе следеће апроксимације:

$$1) \text{MSE}(\hat{t}_{RS}) \approx D(\hat{t}_{RS}) = \sum_h^L \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_h^2(x) + R_h^2 \Delta_h^2(y) - 2R_h \rho_h \Delta_h(x) \Delta_h(y));$$

$\cdot \frac{1}{N^2}$

$$2) \text{MSE}(\bar{x}_{RS}) \approx D(\bar{x}_{RS}) = \frac{1}{N^2} \left[\sum_h^L \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_h^2(x) + R_h^2 \Delta_h^2(y) - 2R_h \rho_h \Delta_h(x) \Delta_h(y)) \right].$$

Доказ: 1) $\text{MSE}(\hat{t}_{RS}) \approx D(\hat{t}_{RS}) = D(\sum_h^L t_h(y) \cdot \hat{R}_h) = \sum_h^L D(t_h(y) \cdot \hat{R}_h) = \sum_h^L D(\hat{t}_{R_h}(x))$
користимо ант. облик за $D(\hat{t}_h)$ из [22] $\approx \sum_h^L \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_h^2(x) + R_h^2 \Delta_h^2(y) - 2R_h \rho_h \Delta_h(x) \Delta_h(y));$
→ стратум = псу ⇒ $\rho(\hat{t}_h)$

2) тривијално (стрелица).

Напомена: За ове дисперзије предлажемо следеће оцене: (Δ_{nh} уместо Δ_h)

$$\widehat{D}(\hat{t}_{RS}) := \sum_h^L \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_{nh}^2(x) + R_h^2 \Delta_{nh}^2(y) - 2R_h \rho_h \Delta_{nh}(x) \Delta_{nh}(y));$$

$\cdot \frac{1}{N^2}$

$$\widehat{D}(\bar{x}_{RS}) := \frac{1}{N^2} \left[\sum_h^L \frac{N_h^2}{n_h} \left(1 - \frac{n_h}{N_h}\right) (\Delta_{nh}^2(x) + R_h^2 \Delta_{nh}^2(y) - 2R_h \rho_h \Delta_{nh}(x) \Delta_{nh}(y)) \right].$$

* Када користимо комбиновану оцену, а када посебну по стратумима?

1° n велико, а нису сви n_h велики: комбинована;

2° n_h велики, а R_h -ови слични: комбинована;

3° n_h велики, а R_h -ови различити: посебна.

* Ако имамо помоћну променљиву u о којој знамо све,
да ли помоћу ње да стратификујемо или да правимо количничку оцену?

Некад једно, некад друго - нема правила.

Најчешће, ако постоји веза која није линеарна \Rightarrow стратификујемо.

Ако веза јесте линеарна \Rightarrow правимо количничку оцену.

\rightarrow због закључка у [22]

Линеарно - регресионо оцењивање

* Као и код количничког, и овде имамо помоћну променљиву

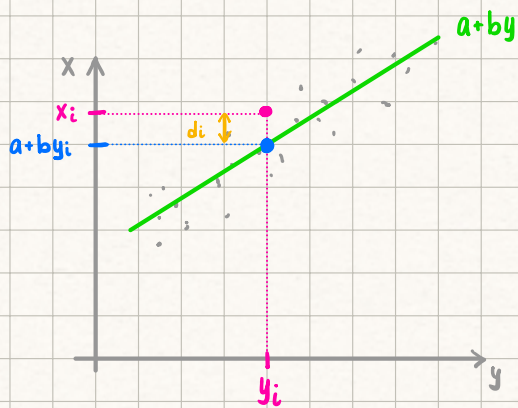
Потребно је: $\rightarrow x, y$ су у корелацији;

\rightarrow из $y=0$ не следи и $x=0$.

Пример: x - принос житарице
 y - загађеност ваздуха

Јесу у корелацији, али ако је $y=0$, принос расте.

* Валимо узорке: $(x_1, y_1), \dots, (x_n, y_n)$.



Узорачки коеф. корелације: $\hat{r} := \frac{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x}_n)(y_i - \bar{y}_n)}{s_n(x) \cdot s_n(y)}$.

Важи: $-1 \leq \hat{r} \leq 1$.

Ако је близу -1 или 1 ,
онда има смисла правити **праву лин. регресије**.

Напомена: $d_i := x_i - (a + by_i)$

27.

Регресионо оцењивање на основу ПСУ без понављања када се b оцењује

Конструишемо праву линеарне регресије $x = a + b \cdot y$, користећи методу најмањих квадрата.

Желимо да збир квадрата одступања (тј. $\sum_{i \in S} d_i^2$) буде што мањи.

Другим речима, тражимо минимум f -је $f(a, b) = \sum_{i \in S} (x_i - (a + b y_i))^2$ (на стд. начин)

$$\frac{\partial f}{\partial a} = 2 \sum (x_i - a - b y_i) \cdot (-1) = 0$$

$$\frac{\partial f}{\partial b} = 2 \sum (x_i - a - b y_i) \cdot (-y_i) = 0$$

$$\sum_{i \in S} x_i - \sum_{i \in S} a - \sum_{i \in S} b y_i = 0$$

$$\Rightarrow n \cdot a = \sum_{i \in S} x_i - b \sum_{i \in S} y_i \quad \stackrel{! : n}{\Rightarrow} \quad a = \bar{x}_n - b \cdot \bar{y}_n$$

$$\sum_{i \in S} x_i y_i - \sum_{i \in S} a y_i - \sum_{i \in S} b y_i^2 = 0$$

$$a \cdot n + b \cdot \sum_{i \in S} y_i = \sum_{i \in S} x_i$$

$$a \cdot \sum_{i \in S} y_i + b \cdot \sum_{i \in S} y_i^2 = \sum_{i \in S} x_i y_i$$

решавамо систем по a, b (преко Крамера: $b = \frac{\Delta_b}{\Delta}$)

$$b = \frac{\begin{vmatrix} n & \sum_{i \in S} x_i \\ \sum_{i \in S} y_i & \sum_{i \in S} x_i y_i \end{vmatrix}}{\begin{vmatrix} n & \sum_{i \in S} y_i \\ \sum_{i \in S} y_i & \sum_{i \in S} y_i^2 \end{vmatrix}} = \frac{n \cdot \sum_{i \in S} x_i y_i - \sum_{i \in S} x_i \cdot \sum_{i \in S} y_i}{n \cdot \sum_{i \in S} y_i^2 - (\sum_{i \in S} y_i)^2} = \frac{\frac{1}{n} \sum_{i \in S} x_i y_i - \bar{x}_n \bar{y}_n}{\frac{1}{n} \sum_{i \in S} y_i^2 - (\bar{y}_n)^2} = \frac{\frac{1}{n} [\sum_{i \in S} x_i y_i - n \cdot \bar{x}_n \bar{y}_n]}{\frac{1}{n} \sum_{i \in S} y_i^2 - (\bar{y}_n)^2}$$

$$\stackrel{(*)}{=} \frac{\frac{1}{n} [\sum_{i \in S} (x_i - \bar{x}_n)(y_i - \bar{y}_n)]}{\stackrel{(**)}{\frac{1}{n} [\sum_{i \in S} (y_i - \bar{y}_n)^2]}} = \frac{(n-1) \hat{\rho} \Delta_n(x) \Delta_n(y)}{(n-1) \Delta_n^2(y)} \quad \Rightarrow \quad b = \hat{\rho} \cdot \frac{\Delta_n(x)}{\Delta_n(y)}$$

$$\begin{aligned} (*) \quad \sum_{i \in S} (x_i - \bar{x}_n)(y_i - \bar{y}_n) &= \sum_{i \in S} x_i y_i - \bar{x}_n \sum_{i \in S} y_i - \bar{y}_n \sum_{i \in S} x_i + \sum_{i \in S} \bar{x}_n \bar{y}_n \\ &= \sum_{i \in S} x_i y_i - n \bar{x}_n \bar{y}_n - n \bar{y}_n \bar{x}_n + n \bar{x}_n \bar{y}_n = \sum_{i \in S} x_i y_i - n \bar{x}_n \bar{y}_n \end{aligned}$$

$$\begin{aligned} (**) \quad \frac{1}{n} \sum_{i \in S} (y_i - \bar{y}_n)^2 &= \frac{1}{n} \sum_{i \in S} y_i^2 - 2 \bar{y}_n \frac{1}{n} \sum_{i \in S} y_i + \frac{1}{n} \sum_{i \in S} \bar{y}_n^2 = \frac{1}{n} \sum_{i \in S} y_i^2 - (\bar{y}_n)^2 \end{aligned}$$

Ово је стационарна тачка, па треба проверити да ли је и минимум.

$$\left. \begin{aligned} \frac{\partial^2 f}{\partial a^2} &= 2n =: A \\ \frac{\partial^2 f}{\partial a \partial b} &= 2 \sum_{i \in S} y_i =: B \\ \frac{\partial^2 f}{\partial b^2} &= 2 \sum_{i \in S} y_i^2 =: C \end{aligned} \right\} \begin{bmatrix} A & B \\ B & C \end{bmatrix} \Rightarrow \begin{aligned} AC - B^2 &> 0 && \text{(екстремум)} \\ A &> 0 && \text{(минимум)} \end{aligned}$$

$$A > 0 \Leftrightarrow 2n > 0, \quad \text{тј.} \quad n > 0 \quad \checkmark$$

$$\begin{aligned} AC - B^2 > 0 &\Leftrightarrow 4n \sum_{i \in S} y_i^2 - 4 \left(\sum_{i \in S} y_i \right)^2 > 0 &\Leftrightarrow n(y_1^2 + \dots + y_n^2) > (y_1 + \dots + y_n)^2 \\ &\Leftrightarrow (n-1)(y_1^2 + \dots + y_n^2) > \underbrace{2y_1 y_2}_{< y_1^2 + y_2^2} + \dots + \underbrace{2y_{n-1} y_n}_{< y_{n-1}^2 + y_n^2} &\Rightarrow \checkmark \\ &&& \text{(сваки се појави } n-1 \text{ пут)} \end{aligned}$$

Лакше, по методу најмањих квадрата, оцене за a и b су: $\hat{b} = \hat{\rho} \cdot \frac{\Delta_n(x)}{\Delta_n(y)}$, $\hat{a} = \bar{x}_n - \hat{b} \cdot \bar{y}_n$.
Тиме добијамо праву лин. регресије: $x = \hat{a} + \hat{b} \cdot y$.
права

Ово све значи да $x_i \approx \hat{a} + \hat{b} y_i$.

Када сумирамо за целу популацију: $\sum x_i \approx N \cdot \hat{a} + \hat{b} \sum y_i \stackrel{:=N}{\Rightarrow} \bar{x} \approx \hat{a} + \hat{b} \bar{y} \stackrel{\hat{a}}{\Rightarrow} \bar{x} \approx \bar{x}_n - \hat{b} \bar{y}_n + \hat{b} \bar{y}$.

Закључак: Ако се бира ПСУ без понављања, при чему b није унапред познат, имамо оцену:

$$\bar{x}_{LR} = \hat{a} + \hat{b} \cdot \bar{y}_n = \bar{x}_n + \hat{b} (\bar{y} - \bar{y}_n)$$

МОТИВ

Теорема 1: Нека се бира ПСУ без понављања. Оцена \bar{x}_{LR} је непристрасна оцена за \bar{x}

Доказ: $E(\bar{x}_{LR}) = E(\bar{x}_n) + \hat{b} (E(\bar{y}) - E(\bar{y}_n)) \stackrel{псу}{=} \bar{x} + \hat{b} (\bar{y} - \bar{y}) = \bar{x};$

Напомена: Може се показати да за велико n : $D(\bar{x}_{LR}) \approx \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot s_x^2 \cdot (1 - \rho^2);$

26.

Регресионо оцењивање на основу ПСУ без понављања када је b_0 познат

деф. Због мотива из [27], за познато b_0 дефинишемо оцене:

$$\bar{X}_{LR} := \bar{X}_n + b_0(\bar{y} - \bar{y}_n);$$

$$\hat{t}_{LR} := N \cdot \bar{X}_{LR}.$$

Теорема 1: Нека се бира ПСУ без понављања и b_0 познато. (претх. искуство)

1) Оцена \bar{X}_{LR} је непристрасна оцена за \bar{X} ;

2) За дисперзију те оцене важи: $D(\bar{X}_{LR}) = \frac{1}{n} \left(1 - \frac{n}{N}\right) (S_x^2 - 2b_0 S_x S_y + b_0^2 S_y^2)$;

3) Непристрасна оцена те дисперзије је $\widehat{D(\bar{X}_{LR})} = \frac{1}{n} \left(1 - \frac{n}{N}\right) (S_n^2(x) - 2b_0 S_n(x) S_n(y) + b_0^2 S_n^2(y))$.

Доказ: 1) $E(\bar{X}_{LR}) = E(\bar{X}_n) + b_0(E(\bar{y}) - E(\bar{y}_n)) \stackrel{НСУ}{=} \bar{X} + b_0(\bar{y} - \bar{y}) = \bar{X}$;

2) $D(\bar{X}_{LR}) = D(\bar{X}_n + b_0(\bar{y} - \bar{y}_n)) \stackrel{(*)}{=} D(\bar{e}_n) \stackrel{НСУ}{=} \left(1 - \frac{n}{N}\right) \frac{1}{n} S_e^2$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i \in S} (e_i - \bar{e})^2 \stackrel{e_i}{=} \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i \in S} [(X_i + b_0(\bar{y} - y_i)) - \bar{X}]^2$$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i \in S} [(X_i - \bar{X}) - b_0(y_i - \bar{y})]^2$$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i \in S} [(X_i - \bar{X})^2 - 2b_0(X_i - \bar{X})(y_i - \bar{y}) + b_0^2(y_i - \bar{y})^2]$$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \left[\underline{S_x^2} - 2b_0 \underline{\hat{\rho}} \underline{S_x S_y} + b_0^2 \underline{S_y^2} \right]$$

(*) $e_i = X_i + b_0(\bar{y} - y_i)$
 $\bar{e} = \bar{X} + b_0(\bar{y} - \bar{y}) = \bar{X}$
 $\bar{e}_n = \bar{X}_n + b_0(\bar{y} - \bar{y}_n)$

3) $\left(1 - \frac{n}{N}\right) \frac{1}{n} S_n^2(e) \stackrel{(**)}{=} \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i \in S} (e_i - \bar{e}_n)^2$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i \in S} [(X_i + b_0(\bar{y} - y_i)) - (\bar{X}_n + b_0(\bar{y} - \bar{y}_n))]^2$$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i \in S} [(X_i - \bar{X}_n) - b_0(\bar{y} - \bar{y}_n - \bar{y} + y_i)]^2$$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i \in S} [(X_i - \bar{X}_n)^2 - 2b_0(X_i - \bar{X}_n)(y_i - \bar{y}_n) + b_0^2(y_i - \bar{y}_n)^2]$$

$$= \left(1 - \frac{n}{N}\right) \frac{1}{n} \left[\underline{S_n^2(x)} - 2b_0 \underline{\hat{\rho}} \underline{S_n(x) S_n(y)} + b_0^2 \underline{S_n^2(y)} \right]$$

$$= \widehat{D(\bar{X}_{LR})}$$

(**) $\left(1 - \frac{n}{N}\right) \frac{1}{n} S_n^2(e)$ је не-пристрасна оцена за $D(\bar{X}_{LR})$ јер је S_n^2 не-пристрасна за S^2 код ПСУ БП. Ми доказујемо да је та оцена исто као наша.

28.

Упоредивање регресионе оцене

* Прво поредимо лин - рег и ПСУ оцене:

$$D(\bar{x}_{LR}) \stackrel{27}{\approx} \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} s_x^2 \underbrace{(1 - \rho^2)}_{-1 \leq \rho \leq 1}, \quad \text{за велико } n$$

$$D(\bar{x}_n) \stackrel{11}{=} \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} s_x^2$$

Закључак: За велико n , \bar{x}_{LR} је боља од \bar{x}_n , осим ако је $\rho = 0$.

* Сада поредимо лин - рег и количничку:

$$D(\bar{x}_R) \stackrel{22}{\approx} \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} (s_x^2 + R^2 s_y^2 - 2R\rho s_x s_y), \quad \text{за велико } n$$

$$\begin{aligned} \Rightarrow D(\bar{x}_R) - D(\bar{x}_{LR}) &= \left(1 - \frac{n}{N}\right) \cdot \frac{1}{n} \cdot (s_x \rho)^2 + (R s_y)^2 - 2R s_y s_x \rho \\ &= \underbrace{\left(1 - \frac{n}{N}\right) \cdot \frac{1}{n}}_{> 0} \cdot \underbrace{(s_x \rho - R s_y)^2}_{\geq 0} \geq 0 \end{aligned}$$

Закључак: За велико n , \bar{x}_{LR} је боља од \bar{x}_R , осим ако је $\rho s_x = R s_y$.

Може се доказати да то важи ако је веза линеарна и пролази кроз коорд. почетак, а већ смо рекли да тада користимо количничку.

Закључак: Линеарно - регресиона оцена је боља и од ПСУ оцене и од количничке оцене.

29.

Регресионо оцењивање на основу стратификованог узорка без понављања комбинована оцена

деф. Комбинована линеарно-регресиона оцена:

$$\bar{X}_{LRC} := \bar{X}_{STR} + b(\bar{y} - \bar{y}_{STR})$$

Теорема 1: Нека се бира стратификован узорак и b_0 познато. (претх. искуство)

1) Оцена \bar{X}_{LRC} је непристрасна оцена за \bar{X} ;

2) За дисперзију те оцене важи:

$$D(\bar{X}_{LRC}) = \frac{1}{N^2} \sum_h^L \frac{1}{n_h} \cdot N_h^2 \left(1 - \frac{n_h}{N_h}\right) \left(S_h^2(x) - 2b_0 \rho_h S_h(x) S_h(y) + b_0^2 S_h^2(y)\right);$$

3) Непристрасна оцена те дисперзије је

$$\widehat{D(\bar{X}_{LRC})} = \frac{1}{N^2} \sum_h^L \frac{1}{n_h} \cdot N_h^2 \left(1 - \frac{n_h}{N_h}\right) \left(S_{n_h}^2(x) - 2b_0 \hat{\rho}_h S_{n_h}(x) S_{n_h}(y) + b_0^2 S_{n_h}^2(y)\right).$$

Доказ: 1) $E(\bar{X}_{LRC}) = E(\bar{X}_{STR}) + b_0(E(\bar{y}) - E(\bar{y}_{STR})) \stackrel{STR}{=} \bar{X} + b_0(\bar{y} - \bar{y}) = \bar{X};$

2) $D(\bar{X}_{LRC}) = D(\bar{X}_{STR} + b_0(\bar{y} - \bar{y}_{STR})) = D\left[\frac{1}{N} \sum_h^L N_h \bar{x}_{n_h} + b_0(\bar{y} - \frac{1}{N} \sum_h^L N_h \bar{y}_{n_h})\right]$

$= D\left[\frac{1}{N} \sum_h^L N_h (\bar{x}_{n_h} + b_0(\bar{y} - \bar{y}_{n_h}))\right] \stackrel{(*)}{=} D\left[\frac{1}{N} \sum_h^L N_h \bar{e}_{n_h}\right]$

независност
стратифика

$= \frac{1}{N^2} \sum_h^L N_h^2 \cdot D[\bar{e}_{n_h}] =$

$\stackrel{PCY}{=} \frac{1}{N^2} \sum_h^L N_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{1}{n_h} \cdot S_h^2(e)$

$= \frac{1}{N^2} \sum_h^L N_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{1}{n_h} \cdot \frac{1}{N_h - 1} \sum_{i=1}^{n_h} (e_{hi} - \bar{e}_h)^2$

$= \frac{1}{N^2} \sum_h^L N_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{1}{n_h} \cdot \frac{1}{N_h - 1} \sum_{i=1}^{n_h} (x_{hi} + b_0(\bar{y} - y_{hi}) - (\bar{x}_h + b_0(\bar{y} - \bar{y}_h)))^2$

$= \frac{1}{N^2} \sum_h^L N_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{1}{n_h} \cdot \frac{1}{N_h - 1} \sum_{i=1}^{n_h} (x_{hi} - \bar{x}_h - b_0(y_{hi} - \bar{y}_h))^2$

$= \frac{1}{N^2} \sum_h^L N_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{1}{n_h} \cdot \frac{1}{N_h - 1} \sum_{i=1}^{n_h} [(x_{hi} - \bar{x}_h)^2 - 2b_0(x_{hi} - \bar{x}_h)(y_{hi} - \bar{y}_h) + b_0^2(y_{hi} - \bar{y}_h)^2]$

$= \frac{1}{N^2} \sum_h^L N_h^2 \cdot \left(1 - \frac{n_h}{N_h}\right) \cdot \frac{1}{n_h} \cdot \left(S_h^2(x) - 2b_0 \rho_h S_h(x) S_h(y) + b_0^2 S_h^2(y)\right);$

(*) $e_{hi} = x_{hi} + b_0(\bar{y} - y_{hi})$
 $\bar{e}_h = \bar{x}_h + b_0(\bar{y} - \bar{y}_h)$
 $\bar{e}_{n_h} = \bar{x}_{n_h} + b_0(\bar{y} - \bar{y}_{n_h})$

$$\begin{aligned}
3) \quad \frac{1}{N^2} \sum_{(h)} N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{1}{n_h} S_{n_h}^2(e) &= \frac{1}{N^2} \sum_{(h)} N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{1}{n_h} \cdot \frac{1}{n_h - 1} \sum_{i \in S} (e_i - \bar{e}_{n_h})^2 \\
&= \frac{1}{N^2} \sum_{(h)} N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{1}{n_h} \frac{1}{n_h - 1} \sum_{i \in S} \left(x_{hi} + b_0(\bar{y} - y_{hi}) - (\bar{x}_{n_h} + b_0(\bar{y} - \bar{y}_{n_h})) \right)^2 \\
&= \frac{1}{N^2} \sum_{(h)} N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{1}{n_h} \frac{1}{n_h - 1} \sum_{i \in S} \left((x_{hi} - \bar{x}_{n_h}) - b_0(\bar{y} - \bar{y}_{n_h} - \bar{y} + y_{hi}) \right)^2 \\
&= \frac{1}{N^2} \sum_{(h)} N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{1}{n_h} \frac{1}{n_h - 1} \sum_{i \in S} \left[(x_{hi} - \bar{x}_{n_h})^2 - 2b_0(x_{hi} - \bar{x}_{n_h})(y_{hi} - \bar{y}_{n_h}) + b_0^2(y_{hi} - \bar{y}_{n_h})^2 \right] \\
&= \frac{1}{N^2} \sum_{(h)} N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{1}{n_h} \left[\hat{S}_{n_h}^2(x) - 2b_0 \hat{\rho}_h \hat{S}_{n_h}(x) \hat{S}_{n_h}(y) + b_0^2 \hat{S}_{n_h}^2(y) \right] \\
&= \hat{D}(X_{LR})
\end{aligned}$$

(**) исти трик
као у [26]

Напомена: Ако b_0 није познато, оцењујемо га са:

$$\hat{b} := \frac{\sum_{(h)} \left[\frac{N_h^2(1 - \frac{n_h}{N_h})}{n_h \cdot (n_h - 1)} \cdot \sum_{i \in S} (x_{hi} - \bar{x}_{n_h})(y_{hi} - \bar{y}_{n_h}) \right]}{\sum_{(h)} \left[\frac{N_h^2(1 - \frac{n_h}{N_h})}{n_h \cdot (n_h - 1)} \cdot \sum_{i \in S} (y_{hi} - \bar{y}_{n_h})^2 \right]}$$

30.

Регресионо оцењивање на основу стратификованог узорка без понављања посебна по стратумима оцена

деф. **Посебна линеарно - регресиона оцена:**

$$\bar{X}_{LRS} := \frac{1}{N} \sum_h^L N_h \cdot \underbrace{(\bar{x}_{n_h} + b_h (\bar{y}_h - \bar{y}_{n_h}))}_{= \bar{x}_{LR}(h)} \quad (\text{стратум} = \text{псу})$$

Теорема 1: Нека се бира стратификован узорак и b_h познати. (претх. искуство)

1) Оцена \bar{X}_{LRS} је непристрасна оцена за \bar{X} ;

2) За дисперзију те оцене важи:

$$D(\bar{X}_{LRS}) = \frac{1}{N^2} \sum_h^L N_h^2 \cdot \frac{1}{n_h} \left(1 - \frac{n_h}{N_h}\right) \left(\Delta_{n_h}^2(x) + b_h^2 \Delta_{n_h}^2(y) - 2b_h \rho_h \Delta_{n_h}(x) \Delta_{n_h}(y) \right);$$

3) Непристрасна оцена те дисперзије је

$$\widehat{D(\bar{X}_{LRS})} = \frac{1}{N^2} \sum_h^L N_h^2 \cdot \frac{1}{n_h} \left(1 - \frac{n_h}{N_h}\right) \left(\Delta_{n_h}^2(x) + b_h \hat{\rho}_h \Delta_{n_h}(x) \Delta_{n_h}(y) \right).$$

Доказ: 1) $E(\bar{X}_{LRS}) = \frac{1}{N} \sum_h^L N_h [E(\bar{x}_{n_h}) + b_h (\bar{y}_h - E(\bar{y}_{n_h}))] \stackrel{псу}{=} \frac{1}{N} \sum_h^L N_h [\bar{x}_h + b_h (\bar{y}_h - \bar{y}_h)] = \frac{1}{N} \sum_h^L N_h \bar{x}_h = \bar{X};$

2) $D(\bar{X}_{LRS}) \stackrel{\text{нез. страт.}}{=} \frac{1}{N^2} \sum_h^L N_h^2 D(\bar{X}_{LR}(h))$
 $\stackrel{26}{=} \frac{1}{N^2} \sum_h^L N_h^2 \cdot \frac{1}{n_h} \left(1 - \frac{n_h}{N_h}\right) \left(\Delta_{n_h}^2(x) + b_h^2 \Delta_{n_h}^2(y) - 2b_h \rho_h \Delta_{n_h}(x) \Delta_{n_h}(y) \right);$


3) $\frac{1}{N^2} \sum_h^L N_h^2 \widehat{D(\bar{X}_{LR}(h))} \stackrel{(**)}{=} \frac{1}{N^2} \sum_h^L N_h^2 \cdot \frac{1}{n_h} \left(1 - \frac{n_h}{N_h}\right) \left(\Delta_{n_h}^2(x) + b_h \hat{\rho}_h \Delta_{n_h}(x) \Delta_{n_h}(y) \right)$
 $= \widehat{D(\bar{X}_{LRS})}$

(**) исти трик као у [26]

Напомена: Ако b_h -ови нису познати, оцењујемо их са: $\hat{b}_h := \hat{\rho}_h \cdot \frac{\Delta_{n_h}(x)}{\Delta_{n_h}(y)}$.

* Када користимо комбиновану оцену, а када посебну по стратумима?

- 1° b_h -ови слични: комбинована;
- 2° b_h -ови различити: посебна.

The background of the page is a grid of light gray lines. Overlaid on this grid is a pattern of hand-drawn teal lines forming a series of overlapping diamond shapes. The diamonds are drawn with varying line thickness and are slightly offset from each other, creating a textured, woven appearance.

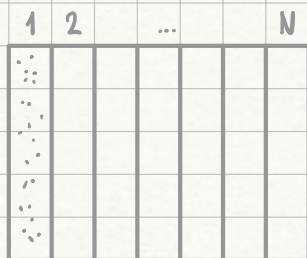
II ЛЕО

НАСТАВАК

Кластер узорак

Идеја је да узорак делимо на дисјунктне подскупове - **кластере**, тј. **примарне јединке**.

Јединке унутар кластера се зову **секундарне јединке**.



Од N кластера,
бирамо њих n .

Напомена: у пракси, бирамо сек. јединке
и онда узмемо цео кластер ком она припада.

Дакле, узоркује се тако што бирамо целе кластере,
тј. све секундарне јединке из кластера које смо изабрали.

N - број примарних јединки;

M_i - број секундарних јединки у i -том кластеру ($i = \overline{1, N}$)

$M_n = \sum_{i=1}^N M_i$ - укупан број секундарних јединки

Кластер узорак користимо када немамо добар узорачки оквир.
Такође, знатно је јефтинији. Ипак, оцене ће бити доста слабије.

Пример: број мобилних телефона у једној кући у Аустралији: кластер = град.

Уместо целу државу да обиђемо, бирамо само градове које обилазимо \Rightarrow јефтино.

Ипак, неки градови су сиромашнији од других \Rightarrow можда слабије оцене.

Кластери треба да буду што хетерогенији унутар себе, а што хомогенији између себе.
Дакле, супротно од стратума.

16.

Кластер узорак

прим. јед. - ПСУ без понављања

| 1 | 2 | 3 | ... | N |
|------------|------------|------------|------------|------------|
| X_{11} | X_{21} | X_{31} | ... | X_{M1} |
| X_{12} | X_{22} | \vdots | ... | X_{M2} |
| \vdots | X_{23} | X_{33} | ... | \vdots |
| X_{1M} | \vdots | | ... | X_{MM} |
| | X_{2M} | | ... | |
| | | | ... | |
| \uparrow | \uparrow | \uparrow | \uparrow | \uparrow |
| M_1 | M_2 | M_3 | | M_M |

деф. $t_i := \sum_j^{M_i} X_{ij}$ - тотална сума i -тог кластера

То значи: $t = \sum_i^M t_i = \sum_i^M \sum_j^{M_i} X_{ij}$.

деф. $\hat{t} := \frac{N}{n} \cdot \left(\sum_{i \in S} t_i \right)$

Теорема 1: 1) Оцена \hat{t} је непристрасна оцена за t ;

2) За дисперзију те оцене важи: $D(\hat{t}) = N^2 \cdot \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_i^M \left(t_i - \frac{t}{N}\right)^2$;

3) Непристрасна оцена те дисперзије је: $\widehat{D}(\hat{t}) := N^2 \cdot \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i \in S} \left(t_i - \frac{1}{n} \sum_{i \in S} t_i\right)^2$.

Доказ: следи из теорије за ПСУ без понављања. (само $x_i \mapsto t_i$)

$$1) E\left(\frac{N}{n} \sum_{i \in S} t_i\right) = N \cdot E\left(\frac{1}{n} \sum_{i \in S} t_i\right) \stackrel{\text{НСУ}}{\stackrel{\text{за } \bar{x}}{=}} N \cdot \frac{1}{N} \sum_i^M t_i = t;$$

$$2) D(\hat{t}) \stackrel{\text{НСУ}}{=} N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} s^2 = N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_i^M \left(t_i - \frac{t}{N}\right)^2;$$

$$3) \widehat{D}(\hat{t}) \stackrel{\text{НСУ}}{=} N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} s_n^2 = N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i \in S} \left(t_i - \frac{\sum_{j \in S} t_j}{n}\right)^2.$$

деф. $\hat{x} := \frac{1}{M_u} \cdot \hat{t}$

Теорема 2: 1) Оцена \hat{x} је непристрасна оцена за \bar{x} ;

2) За дисперзију те оцене важи: $D(\hat{x}) = \frac{1}{M_u^2} \cdot D(\hat{t})$;

3) Непристрасна оцена те дисперзије је: $\widehat{D}(\hat{x}) := \frac{1}{M_u^2} \cdot \widehat{D}(\hat{t})$.

Доказ: следи из Т1.

* Осим ове оцене, постоји и количничка оцена:

Као и до сада, први кластер посматрамо као прву јединку, други као другу итд.
Од N кластера, бирамо n .

Ако изаберемо i -ти кластер, као главно обележје (x) узимамо t_i .
Као помоћно обележје (y), узимамо M_i .

Ово можемо да урадимо: \rightarrow јесу y корелацији
 $\rightarrow y=0 \Rightarrow x=0$

(ако кластер има нула елемената
нема шта да се сабере)

Количник популације је: $R = \frac{\sum_{i=1}^n t_i}{\sum_{i=1}^n M_i} = \frac{t}{M_u}$

Његова оцена је: $\hat{R} = \frac{\sum_{i \in S} t_i}{\sum_{i \in S} M_i}$

Количничке оцене су: 1) $\hat{t}_R = \hat{R} \cdot \sum_{i=1}^n M_i = \hat{R} \cdot M_u$

2) $\bar{x}_R = \frac{1}{M_u} \cdot \hat{t}_R = \hat{R}$

\rightarrow јер су кољ. оцене асимпт. неприс.

Ове оцене су асимптотски непристрасне, па за велико n важи:

$$D(\hat{t}_R) \approx N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i=1}^n (t_i - R M_i)^2$$

Оцена ове дисперзије је:

$$\widehat{D}(\hat{t}_R) = N^2 \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{n-1} \sum_{i=1}^n (t_i - \hat{R} M_i)^2$$

17.

Кластер узорак

прим. јед. - пропорционално

деф. $p_i = \frac{M_i}{M_u}$ је вероватноћа избора i -тог кластера.

Видимо да су вероватноће избора пропорционалне величинама кластера.



Узорак може бити: са понављањем или без понављања.

I) са понављањем - nn оцене.

деф. $\hat{t}_{nn} := \frac{M_u}{n} \cdot \sum_{i \in S} \frac{t_i}{M_i}$.

Теорема 1: 1) Оцена \hat{t}_{nn} је непристрасна оцена за t ;

2) За дисперзију те оцене важи: $D(\hat{t}_{nn}) = \frac{M_u}{n} \cdot \sum_{i \in S} M_i \cdot \left(\frac{t_i}{M_i} - \bar{x} \right)^2$;

3) Непристрасна оцена те дисперзије је: $\widehat{D}(\hat{t}_{nn}) := \frac{1}{n(n-1)} \sum_{i \in S} \left[\frac{M_u}{M_i} \cdot t_i - \hat{t}_{nn} \right]^2$.

Доказ: $p_i = \frac{M_i}{M_u}$ у изразима nn оцена.^[8]

деф. $\bar{x}_{nn} := \frac{1}{M_u} \cdot \hat{t}_{nn}$

Теорема 2: 1) Оцена \bar{x}_{nn} је непристрасна оцена за \bar{x} ;

2) За дисперзију те оцене важи: $D(\bar{x}_{nn}) = \frac{1}{M_u^2} \cdot D(\hat{t}_{nn})$

3) Непристрасна оцена те дисперзије је: $\widehat{D}(\bar{x}_{nn}) := \frac{1}{M_u^2} \cdot \widehat{D}(\hat{t}_{nn})$

Доказ: следи из Т1.

II) Без понављања - немамо адекватну теорију, па користимо НТ оцене (које увек важе)

деф. $\hat{t}_{HT} := \sum \frac{t_i}{\pi_i}$.

Теорема 3: 1) Оцена \hat{t}_{HT} је непристрасна оцена за t ;

2) За дисперзију те оцене важи: $D(\hat{t}_{HT}) = \sum \frac{1-\pi_i}{\pi_i} \cdot t_i^2 + \sum_i \sum_{j \neq i} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} t_i t_j$;

3) Непристрасна оцена те дисперзије је: $\widehat{D}(\hat{t}_{HT}) := \sum \frac{1-\pi_i}{\pi_i^2} \cdot t_i^2 + \sum_i \sum_{j \neq i} \frac{\pi_{ij} - \pi_i \pi_j}{\pi_i \pi_j} \cdot \frac{t_i t_j}{\pi_i \pi_j}$.

Доказ: $p_i = \frac{M_i}{M_u}$ у изразима НТ оцена⁹.

деф. $\bar{x}_{HT} := \frac{1}{M_u} \cdot \hat{t}_{HT}$

Теорема 4: 1) Оцена \bar{x}_{HT} је непристрасна оцена за \bar{x} ;

2) За дисперзију те оцене важи: $D(\bar{x}_{HT}) = \frac{1}{M_u^2} \cdot D(\hat{t}_{HT})$;

3) Непристрасна оцена те дисперзије је: $\widehat{D}(\bar{x}_{HT}) := \frac{1}{M_u^2} \cdot \widehat{D}(\hat{t}_{HT})$.

Доказ: следи из ТЗ.

Напомена: Нама требају π_i и π_{ij} , а знамо само $p_i = \frac{M_i}{M_u}$

За узорковање са понављањем, важи:

1) $\pi_i = 1 - (1-p_i)^n$; (тривијално)

2) $\pi_{ij} = \pi_i + \pi_j - (1 - (1-p_i - p_j)^n)$ (формула укљ. и искљ.: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$)

18.

Упоредивање оцена

ПСУ без понављања и кластер узорак

Претпоставимо да су сви кластери исте величине $M = M_1 = \dots = M_N$.

То значи да је обим популације $N \cdot M$, а обим узорка $n \cdot M$.

деф. **Унутаркластерни коэф. корелације** је коэф. корел. између парова јединки у истом кластеру:

$$r_{uk} = \frac{\sum_i^N \sum_j^M \sum_{k \neq j}^M (x_{ij} - \bar{x})(x_{ik} - \bar{x})}{(M-1)(NM-1)S^2}.$$

Лема 1: 1) $r_{uk} > 0$: кластери су хомогени унутар себе;

2) $r_{uk} < 0$: кластери су хетерогени унутар себе.

Доказ: 1) $r_{uk} > 0 \Rightarrow (x_{ij} - \bar{x})(x_{ik} - \bar{x}) > 0 \Rightarrow$ истог су знака ($++$ или $--$)
 $\Rightarrow x_{ij}, x_{ik}$ су са исте стране \bar{x} $\frac{+}{\bar{x}}$

2) различит знак \Rightarrow са различитих страна \bar{x} .

Лема 2: $D(\hat{t}) = \frac{N^2(1-\frac{n}{N})}{n} \cdot \frac{NM-1}{N-1} \cdot S^2(1+(M-1)r_{uk})$. (алтернативни облик)

Доказ: $D(\hat{t}) \stackrel{(9)(74.2)}{=} N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \sum_i^N (t_i - \frac{t}{N})^2 = N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \sum_i^N (\sum_j^M x_{ij} - M \cdot \frac{t}{N \cdot M})^2$

$$= N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \sum_i^N (\sum_j^M x_{ij} - M \bar{x})^2 = N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \sum_i^N (\sum_j^M (x_{ij} - \bar{x}))^2$$

$$= N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \sum_i^N \left[\sum_j^M (x_{ij} - \bar{x})^2 + \sum_j^M \sum_{k \neq j}^M (x_{ij} - \bar{x})(x_{ik} - \bar{x}) \right]$$

$$= N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \left[\sum_i^N \sum_j^M (x_{ij} - \bar{x})^2 + \sum_i^N \sum_j^M \sum_{k \neq j}^M (x_{ij} - \bar{x})(x_{ik} - \bar{x}) \right]$$

$$= N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \left[(NM-1) \cdot \frac{1}{NM-1} \cdot \sum_i^N \sum_j^M (x_{ij} - \bar{x})^2 + \sum_i^N \sum_j^M \sum_{k \neq j}^M (x_{ij} - \bar{x})(x_{ik} - \bar{x}) \right]$$

$$= N^2(1-\frac{n}{N}) \frac{1}{n} \frac{1}{N-1} \left[(NM-1) \cdot S^2 + (M-1)(NM-1) S^2 r_{uk} \right]$$

$$= N^2(1-\frac{n}{N}) \frac{1}{n} \frac{NM-1}{N-1} \cdot S^2 [1 + (M-1)r_{uk}]$$

Пакле, за велико N (самим тим и велико NM)

$$D(\hat{t}_k) \stackrel{\text{[12]}}{=} \frac{N^2(1-\frac{n}{N})}{n} \cdot \overset{\text{склонили } -1}{\frac{NM}{N}} \cdot S^2(1+(M-1)\rho_{uk}) = N^2(1-\frac{n}{N}) \frac{1}{n} M S^2(1+(M-1)\rho_{uk})$$

$$D(\hat{t}_{psv}) \stackrel{\text{[1]}}{=} N^2 M^2 (1 - \frac{nM}{NM}) \cdot \frac{1}{NM} \cdot S^2 = N^2(1-\frac{n}{N}) \frac{1}{n} M S^2$$

Закључак: $D(\hat{t}_{psv}) < D(\hat{t}_k) \Leftrightarrow \rho_{uk} > 0$. (Тј. по ЛМ: ако су кластери хомогени унутар себе)
лоше

Систематски узорак

деф. Нека је N обим популације, n обим узорка и нека је $k = \frac{N}{n}$.

Од првих k бирамо једну јединку (нпр. i -ту), а онда бирамо сваку k -ту јединку након ње ($i+k, i+2k, \dots$).

Овако добијен узорак зове се **систематски узорак**.

Пример: Ако $\frac{N}{n}$ није цео број, оцене ће бити благо пристрасне

На даље, претпостављамо да јесте цео број, тј. $N = n \cdot k$. Тада су оцене непристрасне.

Приметимо да ово подсећа на стратификован узорак.

Ипак, то није исто - овде имамо неслучајно бирање (остали зависе од првог).

Систематски узорак јесте кластер узорак (имамо k кластера и бирамо један).

↑ први елем.

деф. **Средина систематског узорка** је $\bar{X}_{sis} : \begin{pmatrix} \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_k \\ \frac{1}{k} & \frac{1}{k} & \dots & \frac{1}{k} \end{pmatrix}$. (\bar{x}_i - ср. вр. i -тог могућег сис. узорка)

деф. $S_{sis}^2 := \frac{1}{k(n-1)} \sum_i^k \sum_j^n (x_{ij} - \bar{x}_i)^2$. (x_{ij} је j -ти елемент i -тог узорка)

Теорема 1: 1) Оцена \bar{X}_{sis} је непристрасна оцена за \bar{X} ;

2) За дисперзију те оцене важи: $D(\bar{X}_{sis}) = \frac{N-1}{N} S^2 - \frac{k(n-1)}{N} S_{sis}^2$.

Показ: 1) $E(\bar{X}_{sis}) = \frac{1}{k} \cdot \sum_i \bar{x}_i = \frac{1}{k} \sum_i \frac{1}{n} \sum_j x_{ij} = \frac{1}{kn} \sum_i \sum_j x_{ij} = \frac{1}{N} \cdot t = \bar{X}$;

2) * $D(\bar{X}_{sis}) \stackrel{\text{деф.}}{=} E(\bar{X}_{sis} - \bar{X})^2 \stackrel{\text{деф.}}{=} \frac{1}{k} \cdot \sum_i (\bar{x}_i - \bar{X})^2$

$$* (N-1) S^2 = \sum_i \sum_j (x_{ij} - \bar{X})^2 = \sum_i \sum_j (x_{ij} - \bar{x}_i + \bar{x}_i - \bar{X})^2$$

$$= \sum_i \sum_j (x_{ij} - \bar{x}_i)^2 + \sum_i \sum_j (\bar{x}_i - \bar{X})^2 + 2 \cdot \sum_i (\bar{x}_i - \bar{X}) \cdot \sum_j (x_{ij} - \bar{x}_i)$$

$$\stackrel{(*)}{=} k(n-1) S_{sis}^2 + nk \cdot D(\bar{X}_{sis}) + 0 = k(n-1) S_{sis}^2 + N \cdot D(\bar{X}_{sis})$$

$$\Rightarrow D(\bar{X}_{sis}) = \frac{N-1}{N} S^2 - \frac{k(n-1)}{N} S_{sis}^2$$

$$\begin{aligned} (*) \quad \sum_i (\bar{x}_i - \bar{X}) &= \sum_i \frac{\sum_j x_{ij}}{n} - k \cdot \bar{X} \\ &= \frac{1}{n} \sum_i \sum_j x_{ij} - \frac{N}{n} \cdot \frac{t}{N} \\ &= \frac{t}{n} - \frac{t}{n} = 0 \end{aligned}$$

Последица: \bar{x}_{sis} је боља оцена од \bar{x}_n ако $s^2 < s_{sis}^2$.

Другим речима, \bar{x}_{sis} је боља ако је дисперзија унутар систематских узорака већа, што значи да су систематски узорци (њих k) хетерогенији унутар себе.

- што су хомогенији \Rightarrow дају мању информацију о популацији;
- што су хетерогенији \Rightarrow дају бољу информацију.

Доказ: $D(\bar{x}_{sis}) < D(\bar{x}_n) \Leftrightarrow \frac{N-1}{N} s^2 - \frac{k(n-1)}{N} s_{sis}^2 < (1 - \frac{n}{N}) \frac{1}{n} s^2$

$\Leftrightarrow (\frac{N-1}{N} - \frac{N-n}{Nn}) s^2 < \frac{k(n-1)}{N} s_{sis}^2$

$\Leftrightarrow (1 - \frac{1}{N} - \frac{1}{n} + \frac{1}{N}) s^2 < \frac{k(n-1)}{N} s_{sis}^2$

$\Leftrightarrow \frac{n-1}{n} s^2 < \frac{k(n-1)}{N} s_{sis}^2$

$\Leftrightarrow \stackrel{k=N-n}{s^2} < s_{sis}^2$

* Сада повезујемо ову причу са причом о кластерима.

Теорема 2: $D(\bar{x}_{sis}) = \frac{s^2}{n} \cdot \frac{N-1}{N} \cdot (1 + (n-1) \rho_{uk})$. (алтернативни облик)

Доказ: Као што смо рекли: систематски = кластер узорак где од k кластера дирамо 1.

$\Rightarrow \hat{t}_k = \frac{k}{1} \cdot \frac{\sum_{i \in S} t_i}{k}$ $\bar{x}_{sis} = \frac{\sum_{i \in S} t_i}{n} = \frac{\hat{t}_k}{nk} = \frac{\hat{t}_k}{N}$

↑ др. кластера
↑ само један садржај
↑ др. изабраних

$D(\bar{x}_{sis}) = D(\frac{\hat{t}_k}{N}) = \frac{1}{N^2} D(\hat{t}_k) \stackrel{(16)12}{=} \frac{k^2}{N^2} (1 - \frac{1}{k}) \cdot \frac{1}{1} \cdot \frac{kn-1}{k-1} s^2 (1 + (n-1) \rho_{uk})$

$= \frac{1}{N^2} \cdot \frac{k^2(k-1)}{k} \cdot \frac{kn-1}{k-1} s^2 (1 + (n-1) \rho_{uk}) = \frac{1}{n^2 k^2} \cdot k \cdot (N-1) s^2 (1 + (n-1) \rho_{uk})$

$= \frac{s^2}{n} \cdot \frac{N-1}{N} \cdot (1 + (n-1) \rho_{uk})$

(n) $N \rightarrow k$
M $\rightarrow n$
n $\rightarrow 1$

Напомена: $D(\bar{x}_{sis})$ мање ако је ρ_{uk} мање \Rightarrow кластери што хетерогенији

* Што нисмо транили $\widehat{D}(\bar{x}_{sis})$?

Систематско = кластер за $n=1$.

У кластер оцени $\widehat{D}(\hat{x})$ се појављује члан $\frac{1}{n-1} \Rightarrow$ дељење нулом \downarrow

20.

Квалитет оцене добијене на основу систематског узорка

Први случај

Имамо узорачки списак који је одређен случајно.

(нпр. меримо плату људи, а они су поређани од 1-9999 по последње 4 цифре броја телефона).

Тада ће систематски узорак вероватно дати исти квалитет оцене као и ПСУ.

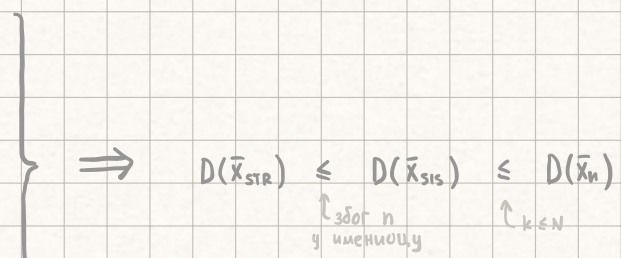
Други случај

Ако је списак направљен и види се да постоји неки линеарни тренд: $x_i = a + b \cdot i$

$$D(\bar{x}_n) = \dots = \frac{(k-1)(N+1)}{12} \cdot b^2; \quad (*)$$

$$D(\bar{x}_{sis}) = \dots = \frac{(k-1)(k+1)}{12} \cdot b^2; \quad (**)$$

$$D(\bar{x}_{str}) = \dots = \frac{(k-1)(k+1)}{12n} \cdot b^2. \quad (***)$$



Дакле, најбољу оцену даје стратификовани узорак.

Трећи случај

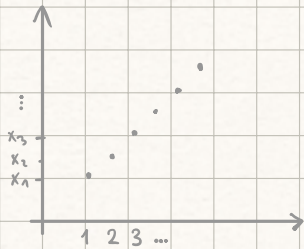
Имамо списак јединки и уочавамо да се вредности обележја периодично понављају.

Има ли смисла користити систематски узорак? Зависи од k .

нпр. 3 1 2 3 1 2 ... и $k=3$: 2 2 2 ... \rightarrow нема смисла;

$k=4$: 2 3 1 2 3 1 \rightarrow букв. идеалан узорак.

Дакле, квалитет оцене зависи од k .



$$x_i = a + b \cdot i, \quad \forall i \Rightarrow \bar{x} = a + b \cdot \frac{N+1}{2}$$

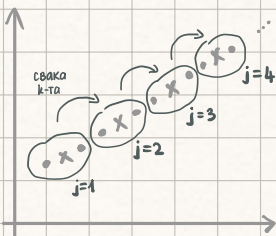
$N = n - k$ - укупан број јединки

n - обим узорка

k - корак код систематског узорка

$$\begin{aligned}
 (*) \quad D(\bar{x}_n) &= \left(1 - \frac{n}{N}\right) \frac{A^2}{n} = \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2 = \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} b^2 \sum_{i=1}^N \left(i - \frac{N+1}{2}\right)^2 \\
 &= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} b^2 \left[\sum_{i=1}^N i^2 - 2 \frac{N+1}{2} \sum_{i=1}^N i + N \cdot \frac{(N+1)^2}{4} \right] \\
 &= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} b^2 \left[\frac{N(N+1)(2N+1)}{6} - \frac{N(N+1)^2}{2} + \frac{N(N+1)^2}{4} \right] \\
 &= \left(1 - \frac{n}{N}\right) \frac{1}{n} \cdot \frac{1}{N-1} b^2 \cdot \frac{N(N+1)}{12} \left[2(2N+1) - 6(N+1) + 3(N+1) \right] = \frac{N-n}{N} \cdot \frac{1}{n} \cdot \frac{1}{N-1} b^2 \cdot \frac{N(N+1)}{12} \cdot (N-1) \\
 &= (nk - n) \frac{1}{n} b^2 \frac{N+1}{12} = \frac{(k-1)(N+1)}{12} b^2
 \end{aligned}$$

(**)



$$j=1: x_{i+0 \cdot k}$$

$$j=2: x_{i+1 \cdot k}$$

$$j=3: x_{i+2 \cdot k}$$

⋮

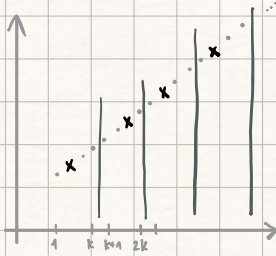
$$x_{i+(j-1) \cdot k}$$

Подсетимо се из 19:

$$\bar{x}_{sis} = \left(\begin{array}{ccc} \bar{x}_1 & \dots & \bar{x}_k \\ 1/k & \dots & 1/k \end{array} \right)$$

$$\begin{aligned}
 D(\bar{x}_{sis}) &= E(\bar{x}_{sis} - \bar{x})^2 = \frac{1}{k} \sum_{j=1}^k (\bar{x}_j - \bar{x})^2 = \frac{1}{k} \sum_{j=1}^k \left[\frac{1}{n} \sum_{i=1}^n x_{i+(j-1) \cdot k} - \left(a + b \frac{N+1}{2}\right) \right]^2 \\
 &= \frac{1}{k} \sum_{j=1}^k \left[a + \frac{bnj}{n} + \frac{bk \sum_{j=1}^n (j-1)}{n} - a - b \frac{N+1}{2} \right]^2 \\
 &= \frac{1}{k} b^2 \sum_{j=1}^k \left[i + \frac{k}{n} \cdot \frac{n(n-1)}{2} - \frac{N+1}{2} \right]^2 = \frac{1}{k} b^2 \sum_{j=1}^k \left[i + \frac{N-k-N-1}{2} \right]^2 = \frac{1}{k} b^2 \sum_{j=1}^k \left(i - \frac{k+1}{2} \right)^2 \\
 &= \frac{1}{k} b^2 \left[\sum_{i=1}^k i^2 - 2 \cdot \frac{k+1}{2} \sum_{i=1}^k i + k \frac{(k+1)^2}{4} \right] = \frac{1}{k} b^2 \left[\frac{k(k+1)(2k+1)}{6} - \frac{k(k+1)^2}{2} + \frac{k(k+1)^2}{4} \right] \\
 &= \frac{1}{k} b^2 k(k+1) \cdot \frac{k-1}{12} = \frac{(k-1)(k+1)}{12} b^2
 \end{aligned}$$

(***)



Како правимо страт. узорак?

Имамо n стратума величине k
и из сваког од њих бирамо узорак обима 1

$$\Rightarrow N_h \rightarrow k, \quad L \rightarrow n, \quad n_h \rightarrow 1$$

Ако бирамо из h -тог стратума, $(h-1) \cdot k$ елемената је прошло и ако узимамо j -ти

$$\Rightarrow X_{(h-1) \cdot k + j}$$

Знамо $D(\bar{X}_{STR}) = \frac{1}{N^2} \sum_h N_h^2 \frac{\Delta_h^2}{n_h} \left(1 - \frac{n_h}{N}\right) \Rightarrow$ морамо још да одредимо $\Delta_h^2 \rightarrow ?$
(све остало познато)

$$\begin{aligned}
 * \Delta_h^2 &= \frac{1}{k-1} \sum_j^k \left[X_{(h-1)k+j} - \frac{1}{k} \sum_l^k X_{(h-1)k+l} \right]^2 \\
 &= \frac{1}{k-1} \sum_j^k \left[a + b((h-1)k+j) - \frac{1}{k} \sum_l^k (a + b((h-1)k+l)) \right]^2 \\
 &= \frac{1}{k-1} \sum_j^k \left[a + b(h-1)k + bj - \frac{1}{k} \cdot ka - \frac{1}{k} \cdot kb(h-1)k + \frac{1}{k} b \sum_l^k l \right]^2 = \frac{1}{k-1} b^2 \sum_j^k \left(j - \frac{k+1}{2} \right)^2 \\
 &= \frac{1}{k-1} b^2 \left[\sum_j^k j^2 - 2 \frac{k+1}{2} \sum_j^k j + k \frac{(k+1)^2}{2} \right] \\
 &= \frac{b^2}{12} k(k+1)
 \end{aligned}$$

$$\begin{aligned}
 * D(\bar{X}_{STR}) &= \frac{1}{N^2} \sum_h N_h^2 \frac{\Delta_h^2}{n_h} \left(1 - \frac{n_h}{N}\right) = \frac{1}{N^2} \sum_h k^2 \cdot \frac{k(k+1)}{12} b^2 \left(1 - \frac{1}{k}\right) \\
 &= \frac{1}{n^2 k^2} \cdot n \cdot k^2 \cdot \frac{k(k+1)}{12} \cdot b^2 \cdot \frac{k-1}{k} \\
 &= \frac{(k^2-1)}{12n} \cdot b^2
 \end{aligned}$$