

# Otanta-aineistojen analyysi

Kevät 2012 Periodi IV

Risto Lehtonen

## Teema 2

### Estimaattoreiden varianssien estimointi

**Estimaattoreiden varianssin estimointi  
linearisointimenetelmällä ja  
pseudotoisto-otannan menetelmillä**

#### **Linearisointimenetelmä**

*Linearization method*

*Taylor series expansion*

#### **Pseudotoisto-otantaan perustuvat menetelmät**

*Pseudoreplication, Sample re-use*

Jackknife-menetelmä (JACKKNIFE)

Balanced Repeated Replications (BRR)

Bootstrap-menetelmä (BOOT)

# **Estimaattoreiden varianssin estimointi linearisointimenetelmällä ja pseudotoisto-otannan menetelmillä**

## **SAS-Ohjelmasovellukset**

### **Yleiskuvaus SAS/STAT 9.2 User's Guide**

#### **SAS/SURVEY-proseduurit**

SURVEYMEANS

SURVEYREG

SURVEYFREQ

SURVEYLOGISTIC

SURVEYPHREG

#### **SAS/STAT Vers. 9.2: Optiot**

##### **Linearisointimenetelmä**

TAYLOR

##### **Pseudotoistomenetelmät**

JACKKNIFE

BRR

**SAS Global Forum 2008**

# LINEARISOINTIMENETELMÄ

Menetelmä on yleisimmin käytetty survey-analyysin ohjelmistoissa:

SAS-proseduurit: SURVEYMEANS,  
SURVEYREG, SURVEYFREQ,  
SURVEYLOGISTIC, SURVEYPHREG  
SPSS:n Complex Surveys-moduli  
Stata: SVY-ohjelmat

## Epälineaariset estimaattorit

Osajoukon koko satunnaismuuttuja  
Osajoukkojen osuusestimaattorit  
Osajoukkojen keskiarvoestimaattorit

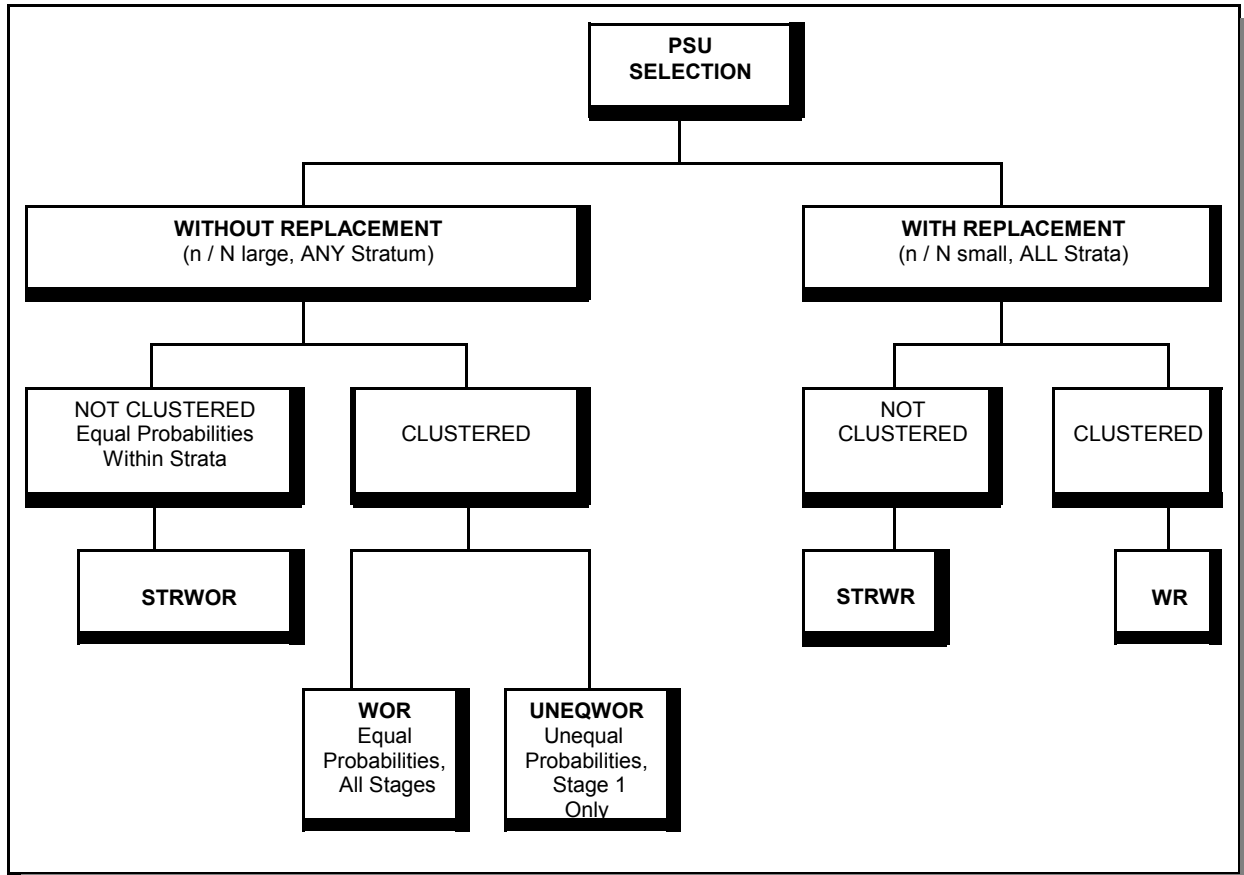
Regressiokertoimien estimaattorit  
Logitmallin kerroinestimaattorit

## HUOM:

**Ohjelmajsovelluksissa (SAS, SPSS, Stata) estimaattoreiden varianssien estimointi perustuu otosrypäiden välisen varianssin estimointiin ositteittain**

**Poikkeus: [SUDAAN](#)-ohjelmisto**

Exhibit 3-1. Choosing the Taylor Series Design Option



# \* PERUSJOUKON OSAJOUKKOJA KOSKEVIEN OSUUKSIEN JA KESKIVARVOJEN ESTIMOINTI

Perusjoukko jaettu  $D$  osajoukkoon  $U_1, \dots, U_D$

**Binäärinen (0/1) indikaattorimuuttuja  $\delta$**

$\delta_{jhik} = 1$  jos ositteen  $h$  rypään  $i$  alkio  $k \in U_j$   
 $= 0$  muulloin

**Binäärinen (0/1) tulosmuuttuja  $y$**

$y_{hik} = 1$  jos ositteen  $h$  rypään  $i$  alkiolla  $k$  on tutkittava ominaisuus  
 $= 0$  muulloin

**Estimoitavana osuusparametri**

$$p_j = \frac{\sum_{h=1}^H \sum_{i=1}^{M_h} \sum_{k=1}^{N_{hi}} \delta_{jhik} y_{hik}}{N_j} = \frac{T_j}{N_j} \quad (j=1, \dots, D)$$

missä

$H$  ositteiden lkm

$M_h$  perusjoukon rypäiden lkm ositteessa  $h$

$N_{hi}$  perusjoukon alkioden lkm ositteen  $h$  rypäessä  $i$

$T_j$  tulosmuuttujan  $y$  totaali osajoukossa  $j$

$N_j$  osajoukon alkioden lkm

**\* SUHTEEN JA OSUUDEN ESTIMAATTORI**  
***Combined ratio estimator***

Osuusestimaattori  $\hat{p}_j$ ,  $j=1, \dots, D$  ( $D$  osajoukkoa)

$$\hat{p}_j = \frac{y_j}{x_j} = \frac{\hat{t}_j}{\hat{N}_j} = \frac{\sum_{h=1}^H \sum_{i=1}^{m_h} y_{jhi}}{\sum_{h=1}^H \sum_{i=1}^{m_h} x_{jhi}} = \frac{\sum_{h=1}^H \sum_{i=1}^{m_h} \sum_{k=1}^{x_{hi}} \delta_{jhik} w_{hik} y_{hik}}{\sum_{h=1}^H \sum_{i=1}^{m_h} \sum_{k=1}^{x_{hi}} w_{hik}}$$

missä  $\hat{t}_j$  Tulosmuuttujan  $y$  totaaliestimaattori osajoukossa  $j$

$\hat{N}_j$  Osajoukon koon estimaattori

$y_j = (n/\hat{N})\hat{t}_j$  ja  $x_j = (n/\hat{N})\hat{N}_j$  vastaavat skaalatut luvut

$w_{hik}$  Painomuuttuja

**HUOM:** Analyttisissä tutkimuksissa painot  $w$  skaalataan usein niin, että niiden keskiarvo koko aineistossa = 1

**Suhteen estimaattori on yksinkertainen esimerkki epälineaarista estimaattorista**

**HUOM:** Merkintätapa  $x_j$  (eikä  $n_j$ ) korostaa sitä, että myös nimittäjä on satunnaismuuttuja

## \* LINEARISOINTIMENETELMÄ

Suhteen estimaattorissa sekä osoittaja  $y_j$  että nimittäjä  $x_j$  ovat satunnaismuuttujia

Tästä syystä asetelmaperusteisen varianssiestimaattorin tulee käsittää:

- osoittajan varianssi  $v(y)$
- nimittäjän varianssi  $v(x)$
- osoittajan ja nimittäjän kovarianssi  $cov(y,x)$

Osajoukon osuustunnusluvun  $\hat{p}_j$  linearisointimenetelmään perustuva varianssiestimaattori on:

$$\hat{V}_{des}(\hat{p}_j) = \hat{p}_j^2 (y_j^{-2} \hat{v}(y_j) + x_j^{-2} \hat{v}(x_j) - 2(y_j x_j)^{-1} cov(y_j, x_j))$$

**HUOM:** Vastaava malliperusteinen (binominen, SRS-perusteinen) varianssiestimaattori:

$$\hat{V}_{bin}(\hat{p}_j) = \hat{p}_j(1 - \hat{p}_j)/\hat{n}_j$$

# LINEARISOINTIMENETELMÄ

[Lehtonen R. and Pahkinen E. \(2004\).](#)

*Practical Methods for Design and Analysis of Complex Surveys. Second Edition. Wiley.*

(materiaali jaetaan luennolla)

Linearisointimenetelmään perustuva  
asetelmaperusteinen

**varianssiapproksimaatio** suhteen  
estimaattorille (osuusestimaattorille)  
(*combined ratio estimator*)

Linearisointimenetelmään perustuva  
osuusestimaattorivektorin

asetelmaperusteinen  
**kovarianssimatriisiestimaattori**

**ESIMERKKI:** OHC Survey

Ositettu kaksiassteinen ryväsotanta



**ESIMERKKI.** Osuusestimaattorin varianssin approksimointi linearisointimenetelmällä  
Lehtonen&Pahkinen 2004, Example 5.5

## OHC Survey demodata

### Ositettu ryväotanta-asetelma

$H= 5$  ositetta

$m= 250$  toimipaikkaa (otosryvästä)

$n = 7841$  henkilöä

### Binäärinen tulosmuuttuja

PHYS Työn fysikaaliset terveyshaitat

0 = Ei ole

1 = On

### Estimointi:

Työn fysikaalisista haitoista kärsivien miesten osuus

### Osuuden estimaatti:

$$\hat{p}_1 = \frac{y_1}{x_1} = \frac{2061}{4485} = 0.4595$$

## Varianssiapproksimaatio:

SAS / SURVEYMEANS

**Osuusestimaattorin asetelmaperusteinen varianssiestimaatti linearisointimenetelmän avulla:**

$$\hat{v}_{des}(\hat{p}_1) = \hat{p}_1^2(y_1^{-2}\hat{v}(y_1) + x_1^{-2}\hat{v}(x_1) - 2(y_1x_1)^{-1}c\hat{ov}(y_1, x_1)) = 0.2775 \times 10^{-3}$$

**SRS-perusteinen (binomimalliin perustuva) varianssiestimaatti:**

$$\hat{v}_{bin}(\hat{p}_1) = \hat{p}_1(1 - \hat{p}_1)/\hat{n}_1 = 0.4595(1 - 0.4595)/4485 = 0.554 \times 10^{-4}$$

**Estimoitu asetelmakerroin:**

$$deff(\hat{p}_1) = 0.0002775/0.0000554 = 5.01$$

Suuri deff-estimaatti viittaa tulosmuuttujan PHYS voimakkaaseen positiiviseen sisäkorrelaatioon rypäissä

Binominen varianssiestimaatti aliestimoii selvästi todellista varianssia

# \* TOISTO-OTANTAAN PERUSTUVA ESTIMAATTORIN VARIANSSIN APPROKSIMOINTI

*Replication / Pseudoreplication methods*

## “Aito” toisto-otanta (*replication*)

a) Perusjoukosta poimitaan useita toisistaan riippumattomia samankokoisia otoksia samalla otanta-asetelmalla niin, että kokonaisotoskoko on  $n$

b) Estimaattoreiden varianssit estimoidaan toisto-otoksista havaitun variaation perusteella

**Käytännössä verraten harvinainen menetelmä**

## “Pseudotoisto”-menetelmät (*pseudoreplication*)

a) Perusjoukosta poimitaan yksi kokoa  $n$  oleva otos annetulla otanta-asetelmalla

b) Poimitusta  $n$  alkion otoksesta poimitaan useita pseudotoisto-otoksia annetulla otanta-asetelmalla

c) Estimaattoreiden varianssit estimoidaan pseudotoisto-otoksista havaitun variaation perusteella

**Käytännössä verraten yleinen menetelmä**

## Teema 2 Estimaattoreiden varianssien estimointi

### Toisto-otanta

(a) Perusjoukosta poimitaan  $K$  toisistaan riippumatonta samankokoista otosta  $s_1, \dots, s_K$  samalla otanta-asetelmalla niin, että kokonaisotoskoko on  $n$

(b) Estimaattoreiden varianssit estimoidaan toisto-otoksista havaitun variaation perusteella

Estimoitava parametri  $\theta = f(T_1, \dots, T_s)$  (totaalien funktio)

Esim:  $\theta = T_1 / T_2$

Estimaattori  $\hat{\theta} = f(\hat{t}_1, \dots, \hat{t}_s)$

Esim:  $\hat{\theta} = \hat{t}_1 / \hat{t}_2$

Otoskoko  $n$

Otokset  $s_1, \dots, s_K$

Otoskoot  $n_1, \dots, n_K, \sum_{k=1}^K n_k = n$

Estimaatit  $\hat{\theta}_1, \dots, \hat{\theta}_K$

Varianssiestimaatti  $v(\hat{\theta}) = \sum_{k=1}^K (\hat{\theta}_k - \hat{\bar{\theta}})^2 / (K - 1)$

$$\hat{\bar{\theta}} = \sum_{k=1}^K \hat{\theta}_k / K$$

## \* PSEUDOTOISTOMENETELMÄT

### Otanta-asetelmat

Perusasetelma: ns. **“Paired clusters design”**

Paljon ositteita

Kustakin ositteesta on poimittu kaksi ryvästä otokseen

Voidaan yleistää mutkikkaampiin asetelmiin, joissa on vaihteleva määrä otosrypäitä per osite

### Estimaattorityypit

Epälineaariset estimaattorit, jotka voidaan lausua totaaliestimaattoreiden funktioina

### Varianssin approksimoinnin perusmenetelmä

Joustava

Soveltuu yleisesti epälineaarille estimaattoreille

Laskentaintensiivinen linearisointimenetelmään verrattuna

## \* PSEUDOTOISTOMENETELMÄT

**Varianssiestimaattorin perusmuoto:**

$$\hat{v}(\hat{\theta}) = c \sum_{k=1}^K (\hat{\theta}_k - \hat{\theta})^2$$

missä  $\hat{\theta}_k$  on pseudo-otoksesta  $k$  laskettu parametrin  $\theta$  estimaatti

$\hat{\theta}$  on alkuperäisestä otoksesta laskettu parametrin  $\theta$  estimaatti

$c$  on vakio, joka riippuu valitusta pseudotoistomenetelmästä

$K$  on kullekin pseudotoistomenetelmälle spesifi toistojen lukumäärä

**HUOM:** Lineaaristen estimaattoreiden tapauksessa kaikki pätevät pseudotoistoperusteiset varianssiestimaatit yhtyvät ja tuottavat vastaavan analyttisen estimaattorin mukaisen estimaatin

**HUOM:** Linearisointimenetelmässä osittaisderivaattojen lausekkeet tarvitaan erikseen kullekin estimaattorityypille

## \* JACKKNIFE-TEKNIikka

### “Jackknife repeated replications” JRR

McCarthy (1966), Frankel (1971), Wolter (1985)

### Pseudo-otosten konstruointi

“Paired clusters design”

$H$  Ositteiden lkm  
 $m_h=2$  Otosrypäitä/osite  
 $n$  Alkiotason otoskoko

### Proseduuri:

1. pseudo-otos:

- a) Poista ensimmäisen ositteen 1. ryväs
- b) Painota toinen ryväs painolla 2
- c) Jätä muut  $H-1$  ositetta ennalleen

Toista proseduuria kullekin  $H$  ositteelle

Saadaan kaikkiaan  $H$  pseudo-otosta (tässä  $K=H$ )

### Komplementtiotokset

Muuta rypäiden poistojärjestys kussakin ositteessa  
Saadaan  $H$  komplementtiotosta

## \* JACKKNIFE-TEKNIikka

### JRR-variانسsiestimaattori “Paired clusters design”

Estimaattoryypit:

Osajoukon osuusestimaattorit

Osajoukon keskiarvoestimaattorit

Regressiokertoimen estimaattorit

Logitmallin kerroinestimaattorit

### JRR-variانسsiestimaattorin perusmuoto:

$$\hat{V}_{JRR}(\hat{\theta}) = \sum_{k=1}^H (\hat{\theta}_k - \hat{\theta})^2$$

**HUOM:** Vakio  $c = 1$  JRR-variانسsiestimaattorin perusmuodolle

Menettelyllä voidaan konstruoida useita vaihtoehtoisia muotoja:

Pseudo-otosten avulla

Komplementtiotosten avulla

Yhdistelmäestimaattoreina

Ks: Lehtonen&Pahkinen (2004) pp. 156-158



## The JRR Technique

The particular jackknife method based on *jackknife repeated replications* has many features of the BRR technique, since only the method of forming the pseudosamples is different. Application of the JRR technique to a design where more than two sample clusters are drawn from a stratum is more straightforward than for BRR. We, however, consider the JRR technique in the simplest case where the number of sample clusters per stratum is exactly two, and the clusters are assumed to be drawn with replacement, i.e. with a design similar to that required for BRR. JRR variance estimators are derived for a ratio estimator  $\hat{r}$ , which is a subpopulation proportion or mean estimator.

We construct the pseudosamples following the method suggested by Frankel (1971). For the first pseudosample, we exclude the first cluster  $h1$  from the first stratum and weight the second cluster  $h2$  by the value 2, leaving the remaining  $H - 1$  strata unchanged. By repeating this procedure for all strata, we get a total of  $H$  pseudosamples. For a similar set of  $H$  complement pseudosamples, we change the order of the clusters that are excluded. The JRR variance estimators are derived using these two sets of pseudosamples.

Like the BRR technique, several alternative JRR variance estimators can be constructed for the parent ratio estimator  $\hat{r}$ . For these, we first derive the pseudosample estimators for each stratum. Let  $\hat{r}_h$  denote a pseudosample estimator based on excluding cluster  $h1$  and duplicating cluster  $h2$  in stratum  $h$ :

$$\hat{r}_h = \frac{2y_{h2} + \sum_{h' \neq h}^H \sum_{i=1}^2 y_{h'i}}{2x_{h2} + \sum_{h' \neq h}^H \sum_{i=1}^2 x_{h'i}}, \quad h = 1, \dots, H. \quad (5.19)$$

These estimators are constructed for each pseudosample. From the complement pseudosamples, we obtain corresponding estimators  $\hat{r}_h^c$  by excluding cluster  $h2$  and duplicating cluster  $h1$ . Using the pseudosample estimators and the complement pseudosample estimators, we can derive the first set of JRR variance estimators for the parent estimator  $\hat{r}$ . Hence we have

$$\hat{v}_{1,jrr}(\hat{r}) = \sum_{h=1}^H (\hat{r}_h - \hat{r})^2, \quad (5.20)$$

and from the complement pseudosamples

$$\hat{v}_{2,jrr}(\hat{r}) = \sum_{h=1}^H (\hat{r}_h^c - \hat{r})^2. \quad (5.21)$$

A combined variance estimator is

$$\hat{v}_{3,jrr}(\hat{r}) = (\hat{v}_{1,jrr}(\hat{r}) + \hat{v}_{2,jrr}(\hat{r}))/2. \tag{5.22}$$

Another set of variance estimators can be obtained using the so-called *pseudovalues* introduced by Quenouille (1956) to reduce the bias of an estimator. In the case considered above, pseudovalues are of the form

$$\hat{r}_h^p = 2\hat{r} - \hat{r}_h, \quad h = 1, \dots, H, \tag{5.23}$$

and for the complement pseudosamples they are denoted by  $\hat{r}_h^{pc}$ . By using the first set of  $H$  pseudovalues  $\hat{r}_h^p$ , we obtain a bias-corrected estimator given by

$$\bar{r}^p = \sum_{h=1}^H \hat{r}_h^p / H, \tag{5.24}$$

and using the pseudovalues  $\hat{r}_h^{pc}$  from the complement pseudosamples we obtain

$$\bar{r}^{pc} = \sum_{h=1}^H \hat{r}_h^{pc} / H. \tag{5.25}$$

Counterparts to the variance estimators (5.20)–(5.22) can be derived from the pseudovalues and the bias-corrected estimators, giving

$$\hat{v}_{4,jrr}(\hat{r}) = \sum_{h=1}^H (\hat{r}_h^p - \bar{r}^p)^2, \tag{5.26}$$

and from the complement pseudosamples

$$\hat{v}_{5,jrr}(\hat{r}) = \sum_{h=1}^H (\hat{r}_h^{pc} - \bar{r}^{pc})^2. \tag{5.27}$$

A combined variance estimator can also be derived:

$$\hat{v}_{6,jrr}(\hat{r}) = (\hat{v}_{4,jrr}(\hat{r}) + \hat{v}_{5,jrr}(\hat{r}))/2. \tag{5.28}$$

Finally, from all the  $2H$  pseudosamples we obtain:

$$\hat{v}_{7,jrr}(\hat{r}) = \sum_{h=1}^H (\hat{r}_h - \hat{r}_h^c)^2 / 4. \tag{5.29}$$

A similar way of constructing the JRR variance estimators was used to that given for the BRR technique. For a linear estimator, the bias-corrected JRR estimators reproduce the parent estimator, and all the JRR variance estimators coincide. This is not the case for nonlinear estimators, but in practice all JRR variance estimators should give closely related results. Like BRR, the variance estimator  $\hat{v}_{7,jrr}$  could be taken as the most natural estimator of the variance of the parent estimator  $\hat{\theta}$ .

The JRR technique can be extended to a more general case in which more than two clusters are drawn from each stratum, for without-replacement sampling of clusters. Pseudosamples and their complements are constructed by consecutively excluding a cluster and weighting the remaining clusters appropriately in a stratum (see Section 4.6 in Wolter 1985).

Like BRR, we use the JRR technique for variance estimation of a ratio estimator  $\hat{r}$  for the MFH Survey design.

### Example 5.3

The JRR technique in the MFH Survey. We continue to consider the estimation of variance of a ratio-type subpopulation proportion estimator  $\hat{p}$  of CHRON (chronic morbidity) and a subpopulation mean estimator  $\bar{y}$  of SYSBP (systolic blood pressure) for 30–64-year-old males. Using the cluster-level data set available, we calculate all the seven JRR variance estimates for  $\hat{p}$  and  $\bar{y}$ .

Because  $H = 24$ , we construct 24 JRR pseudosamples with their complements by the Frankel method. The parent ratio and mean estimates  $\hat{p}$  and  $\bar{y}$ , and the corresponding bias-corrected estimators given by (5.24) and (5.25) based on the pseudovalues  $\hat{p}_h^p, \hat{p}_h^{pc}, \bar{y}_h^p$  and  $\bar{y}_h^{pc}$  calculated from the pseudosamples and their complements, are first obtained. These are

$$\hat{p} = 0.3976, \quad \bar{p}^p = \sum_{k=1}^{24} \hat{p}_k^p / 24 = 0.3972 \quad \text{and} \quad \bar{p}^{pc} = \sum_{k=1}^{24} \hat{p}_k^{pc} / 24 = 0.3980,$$

$$\bar{y} = 141.785, \quad \hat{y}^p = \sum_{k=1}^{24} \bar{y}_k^p / 24 = 141.793 \quad \text{and} \quad \hat{y}^{pc} = \sum_{k=1}^{24} \bar{y}_k^{pc} / 24 = 141.777.$$

All three CHRON proportion estimates and SYSBP mean estimates are close. Next we calculate the JRR variance estimates. For a CHRON proportion estimator  $\hat{p}$  the first variance estimate (5.20) is

$$\hat{v}_{1,jrr}(\hat{p}) = \sum_{h=1}^{24} (\hat{p}_h - 0.3976)^2 = 0.1099 \times 10^{-3},$$

**ESIMERKKI.** Varianssin approksimointi JRR-menetelmällä

## **OHC-demodata**

### **a) Alkuperäinen otanta-asetelma**

Ositettu ryväsotanta-asetelma

$H=5$  ositetta

$m=250$  otosryvästä

$n=7841$  henkilöä

### **b) Modifioitu asetelma**

“Paired clusters design”-asetelma

$H = 125$  ositetta

$m = 250$  otosryvästä

2 otosryvästä per osite

$n = 7841$  henkilöä

## **Binäärinen tulosmuuttuja**

PHYS Työn fysikaaliset terveyshaitat

0 = Ei ole

1 = On

## Estimointi

Työn fysikaalisista haitoista kärsivien miesten osuus

## Osuuden estimaatti

$$\hat{p}_1 = \frac{y_1}{x_1} = \frac{2061}{4485} = 0.4595$$

## Osuuestimaattorin varianssiapproksimaatiot

	Varianssi- estimaatti	<i>deff</i>
--	--------------------------	-------------

---

### a) Alkuperäinen asetelma

JRR	0.0002788	5.03
Linearisointi	0.0002775	5.01

### b) Modifioitu asetelma

JRR	0.0002298	4.15
Linearisointi	0.0002298	4.15

## \* **BOOTSTRAP-TEKNIikka**

### **“Bootstrap repeated replications“**

McCarthy and Snowden (1985)

Rao and Wu (1988)

Rao et al. (1992)

### **Pseudo-otosten konstruointi**

Ositettu ryväotanta-asetelma

$H$  Ositteiden lkm

$m_h = a$  ( $\geq 2$ ) Vakiomäärä otosrypäitä per osite

$n$  Alkiotason otoskoko

Pseudo-otosten konstruointitapa poikkeaa huomattavasti JRR- ja BRR-tekniikoista

Bootstrap on laskentaintensiivisempi tapa

Ei toistaiseksi implementoitu survey-analyysin ohjelmistoihin

## BOOT-proseduuri

**Vaihe 1.** Poimi kokoa  $a$  oleva SRS-WR-otos ositteen  $h$  otosrypäistä,  $h=1, \dots, H$

HUOM: WR-tyyppinen poiminta  
Poiminta suoritetaan toisistaan riippumattomasti jokaisessa  $H$  ositteessa

Saadaan kokoa  $m$  oleva bootstrap-otos

**Vaihe 2.** Toista vaihe 1 kaikkiaan  $K$  kertaa  
(Esim.  $K=1000$ )

Saadaan yhteensä  $K$  riippumatonta bootstrap-otosta

## Bootstrap-varianssiestimaattori

Perusmuoto:

$$V_{BOOT}(\hat{\theta}) = \frac{a}{a-1} \sum_{k=1}^K (\hat{\theta}_k - \hat{\theta})^2 / K$$

**ESIMERKKI:** Tehdään harjoituksissa