

Surveymetodik

Föreläsning 2

Måns Magnusson

Avd. Statistik, LiU

1 Introduktion till urval

- Notation
- Populationsparametrar
- Urvalsteori

2 Obundet slumpmässigt urval

- Inklusionssannolikhet och designvikt
- Estimation vid OSU
- Konfidensintervall vid OSU
- Urvalsdimensionering
- Totaler och proportioner vid OSU

Section 1

Introduktion till urval

N = Antal observationer i populationen

$\mathcal{U} = \{1, 2, 3, \dots, N\}$ = Populationsmängden

n = Antal observationer i urvalet

$\mathcal{S} = \{1, 2, \dots, n\}$ = Urvalsmängden ($\mathcal{S} \subseteq \mathcal{U}$)

i = Observation i

$i \in \mathcal{S}$ = Observation i ingår i urvalsmängden \mathcal{S} , del av urvalet

$\mathcal{K} = \{\mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_K\}$ = Mängden av alla möjliga urval

y = Den egenskap/variabel vi vill undersöka

y_i = egenskapen hos individ i

t = total

p = proportion

\hat{p} = skattning av p

■ Populationstotal

$$t_{\mathcal{U}} = \sum_{i=1}^N y_i = \sum_{i \in \mathcal{U}} y_i$$

■ Populationsmedel

$$\bar{y}_{\mathcal{U}} = \frac{1}{N} \sum_{i=1}^N y_i$$

■ Populationsvarians

$$S_{\mathcal{U}}^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y}_{\mathcal{U}})^2$$

■ Populationsandel

$$p_{\mathcal{U}} = \frac{\text{antal med egenskap av intresse}}{N}$$

- Urvalsteori är en **annan** teori för **inferens** än den “klassiska” statistiska teorin
 - y_i är ett **fixt** värde, **inte** ett utfall av en slumpmässig variabel Y
 - Slumpmässigheten (som vi bygger vår inferens på) **skapar vi själva** genom att dra ett urval slumpmässigt
 - Syftar **bara** till inferens om ändliga populationer (populationsparametrar, inte μ eller σ i en normalfördelning)
 - Brukar ofta kallas **designbaserad** inferens
- Det finns två sätt att dra (skapa) ett slumpmässigt urval
 - med återläggning
 - utan återläggning
 - vidare kommer endast “utan återläggning” behandlas (om inte annat nämns)

Section 2

Obundet slumpmässigt urval

- Obundet slumpmässigt urval - alla urval har samma sannolikhet
- Mest grundläggande urvalsmetod - grunden för andra mer komplicerade metoder
- Antalet tänkbara urval av storlek n från en population av storlek N (utan återläggning):

$$\binom{N}{n} = \frac{N!}{n!(N-n)!}$$

- Urvalsfördelningen för \mathcal{S}_k vid OSU

$$P(\mathcal{S}_k) = \frac{1}{\binom{N}{n}} = \frac{n!(N-n)!}{N!}$$

- Sannolikheten att observation i inkluderas i urvalet: π_i
- **Inklutionssannolikheten** för observation i vid OSU

$$\pi_i = \frac{\# \text{ urval som innehåller obs. } i}{\text{Totalt antal urval}} = (*) = \frac{n}{N}$$

- Summan (över hela populationen) av inklutionssannolikheterna vid OSU är:

$$\sum_{i=1}^N \pi_i = N \cdot \frac{n}{N} = n$$

- **Urvals-/designvikten** (sampling/design weight) för observation i

$$w_i = \frac{1}{\pi_i} = \frac{N}{n}$$

- Summan av vikterna (över urvalet) vid OSU är

$$\sum_{i=1}^n w_i = n \cdot \frac{N}{n} = N$$

- **Tolkning:** Antal personer som representeras av varje respondent i studien

- Vad är sannolikheten för att inkludera den vita bollen i vårt urval om $n = 3$ och $\mathcal{U} = \{\text{Grön, Blå, Röd, Svart, Vit}\}$ ($N = 5$)
- Vi drar ett OSU **utan** återläggning.

- Totalt antal urval vi kan dra är:

$$\binom{N}{n} = \binom{5}{3} = \frac{5!}{(5-3)!3!} = 10$$

- De olika urvalen vi kan dra blir då:

$$\begin{aligned}\mathcal{S}_1 &= \{G, B, R\}, \mathcal{S}_2 = \{G, B, S\}, \mathcal{S}_3 = \{G, B, V\}, \mathcal{S}_4 = \{G, R, S\}, \\ \mathcal{S}_5 &= \{G, R, V\}, \mathcal{S}_6 = \{G, S, V\}, \mathcal{S}_7 = \{B, R, S\}, \mathcal{S}_8 = \{B, R, V\}, \\ \mathcal{S}_9 &= \{B, S, V\}, \mathcal{S}_{10} = \{R, S, V\}\end{aligned}$$

- Sannolikheten för respektive urval är $P(\mathcal{S}_k) = \frac{1}{10}$, vilket ger att

$$\pi_{\text{Vit}} = P(\mathcal{S}_3) + P(\mathcal{S}_5) + P(\mathcal{S}_6) + P(\mathcal{S}_8) + P(\mathcal{S}_9) + P(\mathcal{S}_{10}) = \frac{3}{5}$$

- Eller vi kan använda formeln på slide 10:

$$\pi_{\text{Vit}} = \frac{n}{N} = \frac{3}{5}$$

Subsection 2

Estimation vid OSU

- Vi är intresserade av att skatta eller **estimera** populationsparametrarna $\bar{y}_{\mathcal{U}}$ och $S_{\mathcal{U}}^2$.
- För detta använder vi **estimatorer**.

$$\hat{y}_{\mathcal{U}} = \bar{y}_{\mathcal{S}} = \frac{1}{n} \sum_{i \in \mathcal{S}} y_i$$

$$\hat{S}_{\mathcal{U}}^2 = S_{\mathcal{S}}^2 = s^2 = \frac{1}{n-1} \sum_{i \in \mathcal{S}} (y_i - \bar{y})^2$$

- På grund av urvalfördelningen $P(\mathcal{S})$ får estimatorerna en sannolikhetsfördelning.

- Vi fortsätter på exemplet på slide 12 ovan där $N = 5$, $n = 3$ och $\mathcal{U} = \{\text{Grön, Blå, Röd, Svart, Vit}\}$ och $P(\mathcal{S}) = \frac{1}{10}$
- De olika bollarna väger olika mycket

$$y = \{5, 10, 8, 7, 12\}$$

- Populationsparametrar (de sanna värdena) som vi är intresserade av är:
- Populationsmedelvärdet:

$$\bar{y}_U = \frac{1}{N} \sum_{i \in U} y_i = \frac{1}{5} \sum_{i \in \{a,b,c,d,e\}} y_i = 8.4$$

- Populationsvarians/standardavvikelse:

$$S_U^2 = \frac{1}{N-1} \sum_{i \in U} (y_i - \bar{y}_U)^2 = \frac{1}{4} \sum_{i \in \{a,b,c,d,e\}} (y_i - 8.4)^2 = 7.3$$

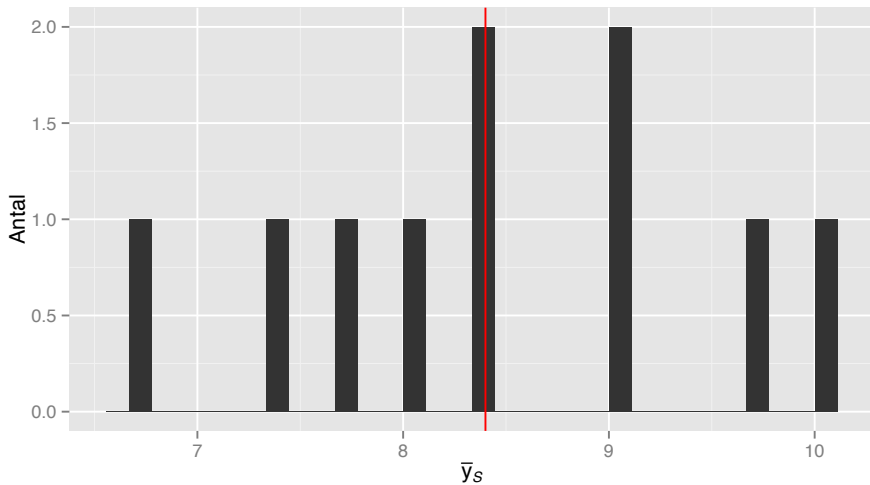
$$S_U = \sqrt{7.3} \approx 2.3$$

- Vad är samplingfördelningen för estimatorn \hat{y}_U ?

- Med R kan vi beräkna samplingfördelningen för \bar{y}_S då $n = 3$ och $y = (5, 10, 8, 7, 12)$
- Här används funktionen `samplingDist()` som finns [här](#).

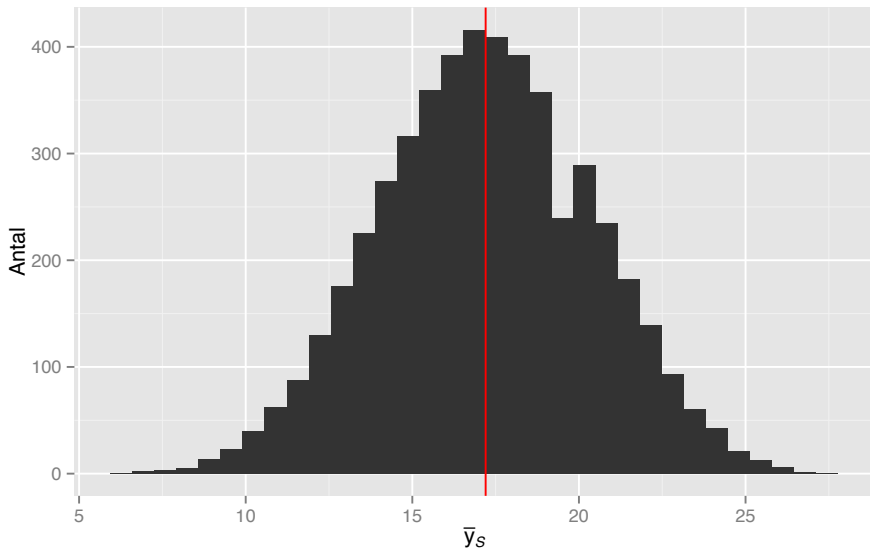
| | Obs.1 | Obs.2 | Obs.3 | P_S | y_hat | s_hat |
|----|-------|-------|-------|-----|-------|-------|
| 1 | 5 | 10 | 8 | 0.1 | 7.67 | 2.52 |
| 2 | 5 | 10 | 7 | 0.1 | 7.33 | 2.52 |
| 3 | 5 | 10 | 12 | 0.1 | 9.00 | 3.61 |
| 4 | 5 | 8 | 7 | 0.1 | 6.67 | 1.53 |
| 5 | 5 | 8 | 12 | 0.1 | 8.33 | 3.51 |
| 6 | 5 | 7 | 12 | 0.1 | 8.00 | 3.61 |
| 7 | 10 | 8 | 7 | 0.1 | 8.33 | 1.53 |
| 8 | 10 | 8 | 12 | 0.1 | 10.00 | 2.00 |
| 9 | 10 | 7 | 12 | 0.1 | 9.67 | 2.52 |
| 10 | 8 | 7 | 12 | 0.1 | 9.00 | 2.65 |

Samplingfördelningen för \bar{y}_S

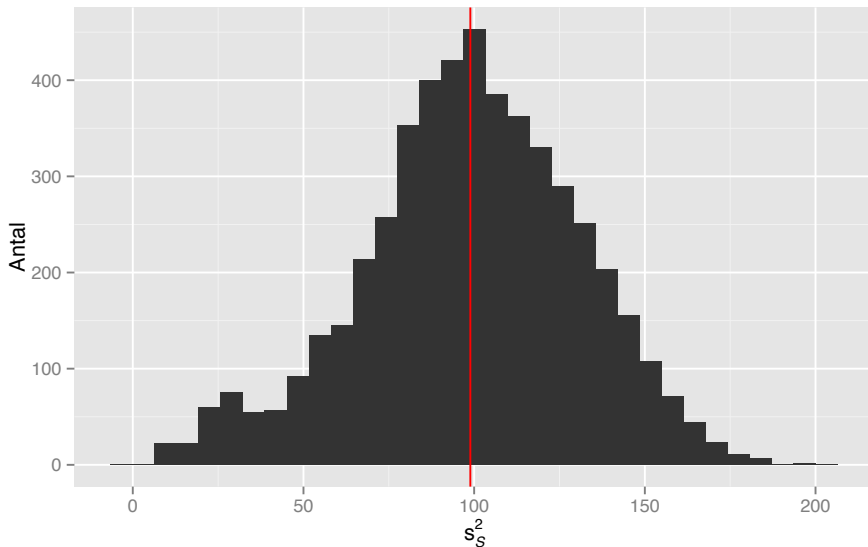


- Samplingfördelningen för \bar{y}_S då $n=6$, $N=15$ och $\binom{15}{6} = 5005$.
- $y_U = (29, 21, 23, 3, 22, 30, 24, 6, 15, 2, 4, 10, 16, 27, 26)$.

| | Obs.1 | Obs.2 | Obs.3 | Obs.4 | Obs.5 | Obs.6 | P_S | y_hat | s_hat |
|----|-------|-------|-------|-------|-------|-------|--------|-------|-------|
| 1 | 29 | 21 | 23 | 3 | 22 | 30 | 0.0002 | 21.3 | 9.73 |
| 2 | 29 | 21 | 23 | 3 | 22 | 24 | 0.0002 | 20.3 | 8.94 |
| 3 | 29 | 21 | 23 | 3 | 22 | 6 | 0.0002 | 17.3 | 10.37 |
| 4 | 29 | 21 | 23 | 3 | 22 | 15 | 0.0002 | 18.8 | 8.95 |
| 5 | 29 | 21 | 23 | 3 | 22 | 2 | 0.0002 | 16.7 | 11.33 |
| 6 | 29 | 21 | 23 | 3 | 22 | 4 | 0.0002 | 17.0 | 10.83 |
| 7 | 29 | 21 | 23 | 3 | 22 | 10 | 0.0002 | 18.0 | 9.59 |
| 8 | 29 | 21 | 23 | 3 | 22 | 16 | 0.0002 | 19.0 | 8.88 |
| 9 | 29 | 21 | 23 | 3 | 22 | 27 | 0.0002 | 20.8 | 9.26 |
| 10 | 29 | 21 | 23 | 3 | 22 | 26 | 0.0002 | 20.7 | 9.14 |
| 11 | 29 | 21 | 23 | 3 | 30 | 24 | 0.0002 | 21.7 | 9.79 |
| 12 | 29 | 21 | 23 | 3 | 30 | 6 | 0.0002 | 18.7 | 11.54 |
| 13 | 29 | 21 | 23 | 3 | 30 | 15 | 0.0002 | 20.2 | 10.05 |
| 14 | 29 | 21 | 23 | 3 | 30 | 2 | 0.0002 | 18.0 | 12.49 |



■ Samplingfördelningen för s_S^2



- Oftast är vi intresserade av **osäkerheten** i de punktskattningar vi producerar.
- Osäkerheten kommer från urvalsfördelningen $P(S)$
- Om vi känner till $S_{\mathcal{U}}^2$ (vilket vi aldrig gör) är

$$\text{Var}(\hat{y}_{\mathcal{U}}) = \frac{S_{\mathcal{U}}^2}{n} \left(1 - \frac{n}{N}\right)$$

- Istället använder vi punktskattningen för $S_{\mathcal{U}}^2$, nämligen s^2 och **uppskattar** vår osäkerhet.

$$\hat{\text{Var}}(\hat{y}_{\mathcal{U}}) = \frac{s^2}{n} \left(1 - \frac{n}{N}\right)$$

- $\left(1 - \frac{n}{N}\right)$ kallas **ändlighetskorrektion** (finite population corrections)
- $\frac{n}{N}$ kallas (i dessa sammanhang) **urvalsfraction** (sampling fraction)

- Istället för en punktskattnings varians brukar ofta en punktskattnings **medelfel** (standardfel, standard error) användas

$$\hat{SE}(\hat{y}_U) = \sqrt{\hat{Var}(\hat{y}_U)} = \sqrt{\frac{s^2}{n} \left(1 - \frac{n}{N}\right)}$$

- Medelfel $\hat{SE}(\hat{y}_U)$ och **urvalsstandardavvikelsen** s^2 är olika saker.

Subsection 3

Konfidensintervall vid OSU

- För att **beskriva** vår osäkerhet bildar vi ofta **konfidsensintervall** (KI) för våra estimatorer $\hat{y}_{\mathcal{U}}$
- Inom urvalsteorin antar vi **inte** någon fördelning/modell för y - vi har bara samplingfördelningen $P(\mathcal{S})$
- Normalapproximation kan göras med **centrala gränsvärdessatsen (CGS)** om:
 - N och n är "stora" och
 - Fördelningen i urvalet, $y_{\mathcal{S}}$ inte är "alltför" skev
 - **Varning** för variabler med (vilka ofta kan vara skeva):
 - små proportioner
 - inkomster
- För tumregler, se Lohr (2009, formel 2.23, s. 44)

- Om \bar{y}_S är approximativt normalfördelat så

$$\frac{\hat{y}_U - \bar{y}_U}{SE(\bar{y}_S)} \sim N(0, 1)$$

- Så ett $100 \cdot (1 - \alpha)\%$ KI för \bar{y}_U är

$$\bar{y}_S \pm z_{\alpha/2} \cdot SE(\bar{y}_S) = \bar{y}_S \pm z_{\alpha/2} \cdot \sqrt{\frac{S_U^2}{n} \left(1 - \frac{n}{N}\right)}$$

där $z_{\alpha/2}$ är den $\alpha/2$ percentilen i normalfördelningen

- $z_{\alpha/2} \cdot \hat{SE}(\bar{y}_S)$ brukar kallas för **felmarginal**.
- t -fördelningen används ofta i statistikprogram som R och SAS.
 - Detta för att ta hänsyn till osäkerheten i s^2

- Konfidensintervall för \hat{y}_U då $n = 6$, $N = 15$ från det föregående exemplet.
- Det sanna värdet i populationen: $\bar{y}_U = 17.2$.
- Nedan följer de 8 första möjliga urvalen (kombinationer) av totalt $K = 5005$ teoretiskt möjliga urval.

| | Obs..1 | Obs..6 | P_S | y_hat | s_hat | SE_hat | t | KI.low | KI.up | in.KI |
|---|--------|--------|--------|-------|-------|--------|------|--------|-------|-------|
| 1 | 29 | 30 | 0.0002 | 21.3 | 9.73 | 3.08 | 2.57 | 13.42 | 29.2 | 1 |
| 2 | 29 | 24 | 0.0002 | 20.3 | 8.94 | 2.83 | 2.57 | 13.07 | 27.6 | 1 |
| 3 | 29 | 6 | 0.0002 | 17.3 | 10.37 | 3.28 | 2.57 | 8.91 | 25.8 | 1 |
| 4 | 29 | 15 | 0.0002 | 18.8 | 8.95 | 2.83 | 2.57 | 11.56 | 26.1 | 1 |
| 5 | 29 | 2 | 0.0002 | 16.7 | 11.33 | 3.58 | 2.57 | 7.46 | 25.9 | 1 |
| 6 | 29 | 4 | 0.0002 | 17.0 | 10.83 | 3.42 | 2.57 | 8.20 | 25.8 | 1 |
| 7 | 29 | 10 | 0.0002 | 18.0 | 9.59 | 3.03 | 2.57 | 10.20 | 25.8 | 1 |
| 8 | 29 | 16 | 0.0002 | 19.0 | 8.88 | 2.81 | 2.57 | 11.78 | 26.2 | 1 |

- Täckningsgraden (baserat på t -fördelningen) för konfidensintervallen är 0.946.

Subsection 4

Urvalsdimensionering

- Hur stort urval ska vi dra? Vanlig (kostnads-) fråga!
- Beror på hur säkra vi vill vara. Bredd på KI.
- Vad kan vi påverka i formeln för KI?

$$\bar{y}_S \pm z_{\alpha/2} \cdot \hat{SE}(\bar{y}_S) = \bar{y}_S \pm z_{\alpha/2} \cdot \sqrt{\frac{s^2}{n} \left(1 - \frac{n}{N}\right)}$$

- Minska konfidensgraden (α)
 - ett KI med $\alpha = 0.2$ är kortare än ett med $\alpha = 0.01$
- Urvalsstorleken (n)
 - ett större urval ger ett smalare KI

- Hur stort urval ska vi dra för att få en viss bredd på KI?
 - Vi behöver veta $S_{\mathcal{U}}^2$
(tar vi från annan undersökning eller gissar)
- Låt e vara den felmarginal vi vill ha

$$e = z_{\alpha/2} \cdot \sqrt{\frac{S_{\mathcal{U}}^2}{n} \left(1 - \frac{n}{N}\right)}$$

- Lös ut n så fås (*)

$$n = \frac{n_0}{1 + \frac{n_0}{N}}, \text{ där } n_0 = \frac{z_{\alpha/2}^2 S_{\mathcal{U}}^2}{e^2}$$

- Då $N \rightarrow \infty$ så går $n \rightarrow n_0$

- Vi vet vill dimensionera en studie för att uppskatta kostnaderna för magsjuka i Katrineholm med $N = 4328$.
- Baserat på tidigare studier antar vi att $S_U = 938$.
- Vi är intresserade av ett konfidensintervall (95 %) på maximalt ± 100 kr.

Vi har att $N = 4328$, $z_{2.5\%} = 1.96$, $e = 100$ och $S_U = 938$.

Vi använder oss av:

$$n_0 = \frac{z_{\alpha/2}^2 S_U^2}{e^2} = \frac{1.96^2 \cdot 938^2}{100^2} = 338.001$$

Sedan använder vi

$$n = \frac{n_0}{1 + \frac{n_0}{N}} = \frac{338.001}{1 + \frac{338.001}{4328}} = 313.516 \approx 314$$

Således behöver vi i detta fall ett urval på 314.

Subsection 5

Totaler och proportioner vid OSU

- Vanliga skattningar av intresse är **totaler** och **proportioner**
- Båda är specialfall av \bar{y}
- Total (t)
 - Populationsparameter:

$$t_{\mathcal{U}} = \sum_{i \in \mathcal{U}} y_i = \sum_{i=1}^N y_i = N\bar{y}_{\mathcal{U}}$$

- Proportion (p)
 - Populationsparameter av intresse:

$$p_{\mathcal{U}} = \frac{\# \text{ med egenskap av intresse}}{N}$$

- Om y antar värden $\{0, 1\}$ så

$$p_{\mathcal{U}} = \frac{1}{N} \sum_{i=1}^N y_i = \bar{y}_{\mathcal{U}}$$

- På följande sätt estimeras totalen (**totalestimator**) vid OSU

$$\hat{t}_{\mathcal{U}} = N\bar{y}_{\mathcal{S}} = \frac{N}{n} \sum_{i \in \mathcal{S}} y_i$$

- Med medelfelet (*)

$$\hat{SE}(\hat{t}_{\mathcal{U}}) = N \cdot \sqrt{\frac{s_y^2}{n} \left(1 - \frac{n}{N}\right)}$$

- Detta ger följande KI för $\hat{t}_{\mathcal{U}}$

$$N\bar{y}_{\mathcal{S}} \pm z_{\alpha/2} \cdot N \cdot \sqrt{\frac{s_y^2}{n} \left(1 - \frac{n}{N}\right)}$$

Obs! s_y^2 är den uppskattade populationsvariansen för y .

- Vi vill uppskatta de totala sjukvårdskostnaderna för magsjuka i befolkningen (15 - 64 år).
- Vi har gjort en undersökning med

$$n = 2731$$

$$N = 4942059$$

$$\bar{y}_S = 76.5$$

$$s = 191$$

- Beräkna en skattning av de totala sjukvårdskostnaderna för magsjuka med tillhörande 95%-igt konfidensintervall.

- Detta ger följande totalskattning:

$$\hat{t}_{\mathcal{U}} = N\bar{y}_S = 4942059 \cdot 76.5 = 378067513.5$$

med konfidensintervallet

$$\begin{aligned} & \hat{t}_{\mathcal{U}} \pm z_{\alpha/2} \cdot N \cdot \sqrt{\frac{s_y^2}{n} \left(1 - \frac{n}{N}\right)} \\ = & 378067513.5 \pm 1.96 \cdot 4942059 \cdot \sqrt{\frac{191^2}{2731} \left(1 - \frac{2731}{4942059}\right)} \\ = & 378067513.5 \pm 1.96 \cdot 4942059 \cdot \sqrt{\frac{191^2}{2731} (0.999)} \\ = & 378067513.5 \pm 1.96 \cdot 18057616.004 \\ \Rightarrow & [413460440.868 : 342674586.132] \end{aligned}$$

- På följande sätt estimeras proportionen (då $y = \{0, 1\}$)

$$\hat{p}_U = p_S = \bar{y}_S = \frac{1}{n} \sum_{i \in S} y_i$$

- Standardfelet (*)

$$\hat{SE}(\hat{p}) = \sqrt{\frac{\hat{p}(1 - \hat{p})}{n - 1} \left(1 - \frac{n}{N}\right)}$$

är ett specialfall då $y = \{0, 1\}$.

- Vilket ger följande KI för \hat{p}_U

$$\hat{p}_U \pm z_{\alpha/2, n-1} \cdot \sqrt{\frac{\hat{p}(1 - \hat{p})}{n - 1} \left(1 - \frac{n}{N}\right)}$$

- Vi vill estimera andelen personer i befolkningen 16-79 år som utsatts för misshandel under 2011. (se Irlander and Hvitfeldt (2012))

Slog, sparkade eller utsatte någon dig med avsikt för något annat fysiskt våld, så att du skadades eller så att det gjorde ont, under förra året (2011)?

- Följande uppgifter har vi

$$N = 7297354$$

$$n = 13386$$

$$\sum_{i \in \mathcal{S}} y_i = 363$$

- **Obs!** I NTU används stratifierat urval och kalibrerade vikter, så “på riktigt” får vi en annan siffra.

- Detta ger följande skattning av populationsandelen:

$$\hat{p}_U = \frac{1}{n} \sum_{i \in S} y_i = \frac{363}{13386} = 0.027$$

med konfidensintervallet

$$\begin{aligned} & \hat{p}_U \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{p}(1 - \hat{p})}{n - 1} \left(1 - \frac{n}{N}\right)} \\ = & 0.027 \pm 1.96 \cdot \sqrt{\frac{0.027(1 - 0.027)}{13386 - 1} \left(1 - \frac{13386}{7297354}\right)} \\ = & 0.027 \pm 1.96 \cdot \sqrt{\frac{0.026}{13385} (0.998)} \\ = & 0.027 \pm 1.96 \cdot 0.001 \\ \Rightarrow & [0.03 : 0.024] \end{aligned}$$

- Brås skattning är 0.025.

- Irlander, Å., Hvitfeldt, T., 2012. NTU 2011 : om utsatthet, trygghet och förtroende. Brottsförebyggande rådet, Stockholm.
- Lohr, S., 2009. Sampling: design and analysis, 2nd Edition. Thomson.