

Census like population size estimation based on administrative data ---

Some recent developments on capture
recapture methodology in the presence of
erroneous enumerations

Li-Chun Zhang

(L.Zhang@soton.ac.uk)

University of Southampton & Statistics Norway

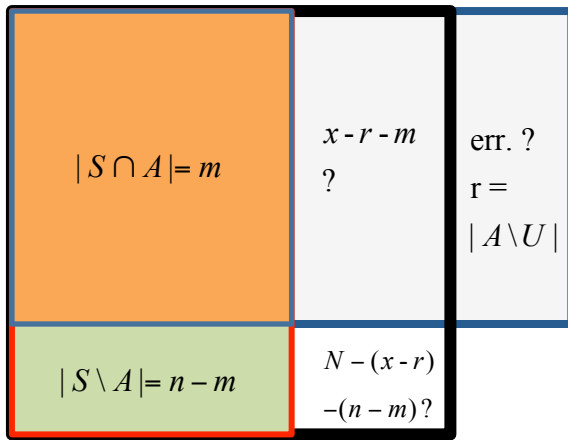
John Dunne

Central Statistics Office, Ireland, John.Dunne@cso.ie

To start with, a traditional census costs too much, takes too long...

- **Alternative**: population register (**PR**) for some only, **administrative registers for all**
- **Q1**: What if one replaces the census with a register?
- Trimmed Dual System Estimation (TDSE)
 - Basic idea; additional complications; possibly combining with additional modelling
- Potential areas of application
 - Population size estimation: regular or irregular
 - Health screening , edit scoring, etc.
- The Irish question (**Q2**): Can one estimate population size **without census, nor PR, nor survey**?

Survey S
(size n)



List A
(size x)

Target population U
(size N, unknown)

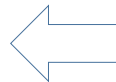


$$\tilde{N} = n \frac{x}{m} \geq$$

Naive DSE

$$\tilde{N} = n \frac{x - r}{m}$$

Hypothetic DSE



Assume:

$$P(I_{iS} = 0 | I_{iU} = 1) = P\left(I_{iS} = 0 \mid \begin{matrix} I_{iU} = 1 \\ I_{iA} = 1 \end{matrix}\right) = P\left(I_{iS} = 0 \mid \begin{matrix} I_{iU} = 1 \\ I_{iA} = 0 \end{matrix}\right)$$

$$P(I_{iS} = 1 | I_{iU} = 1) = P\left(I_{iS} = 1 \mid \begin{matrix} I_{iU} = 1 \\ I_{iA} = 1 \end{matrix}\right) = P\left(I_{iS} = 1 \mid \begin{matrix} I_{iU} = 1 \\ I_{iA} = 0 \end{matrix}\right)$$

$$\frac{N - E(n)}{E(n)} = \frac{E(x - r - m)}{E(m)} = \frac{E[N - (x - r) - (n - m)]}{E(n - m)}$$

$$\frac{N}{E(n)} = \frac{E(x - r)}{E(m)}$$

$$N = E(n) \frac{E(x - r)}{E(m)}$$

Trimmed DSE (TDSE)

k_1 units in S and $A + k_0$ units in $A \setminus S$

Ideal scoring
(k, k_1) = ($r, 0$)

Scoring k units in A

No scoring $k=0$
Naïve DSE

$$\tilde{N} = n \frac{x - r}{m - 0}$$

TDSE _{k}

$$\hat{N}_0 = n \frac{x}{m}$$

$$\hat{N}_k = n \frac{x - k}{m - k_1}$$

$\tilde{N} \leq \hat{N}_k$ iff

$$\frac{x - r}{m} \leq \frac{x - k}{m - k_1} \Leftrightarrow \frac{k - r}{x - r} \leq \frac{k_1}{m}$$

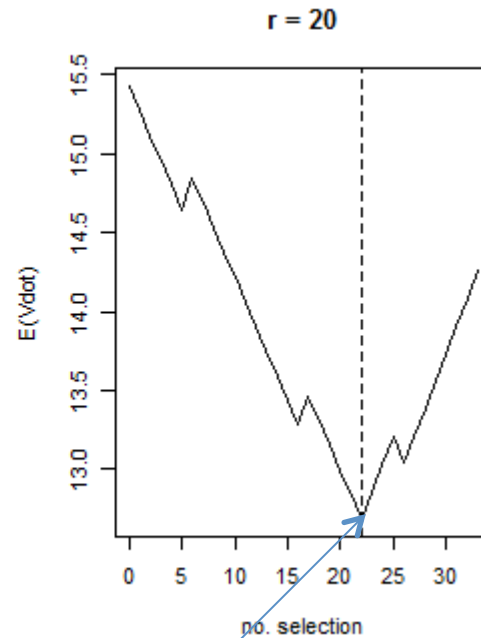
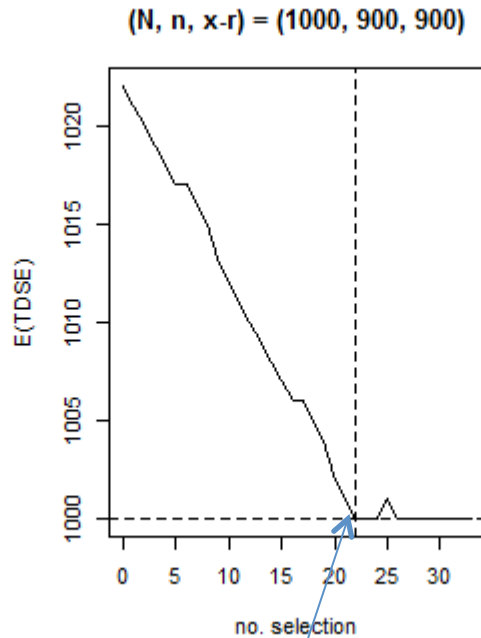
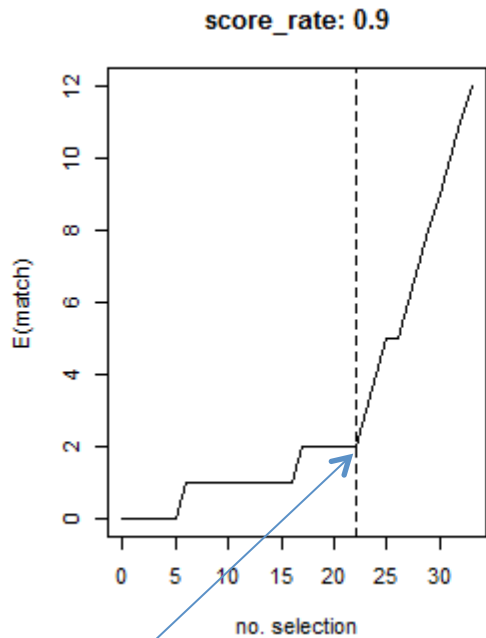
$\hat{N}_k \leq \hat{N}_0$ iff

$$\frac{x - k}{m - k_1} \leq \frac{x}{m} \Leftrightarrow \frac{k_1}{m} \leq \frac{k}{x}$$

OK as long as $k < r$
i.e. "safe" up to $k = r$

Needs to be more effective than random scoring

Stopping rule



$$k = k_1 + k_{0u} + k_{0r}$$

$$p_{rk} = E(k_{0r}) / k \Rightarrow k_r = r / p_{rk}$$

$$\frac{E(k_1)}{k} = \begin{cases} \left[(1 - p_{rk}) \frac{m}{x - r} \right] & \text{if } k \leq k_r \\ \left[(1 - \frac{r}{k}) \frac{m}{x - r} \right] & \text{if } k > k_r \end{cases}$$

$$E(\hat{N}_k) = E\left(n \frac{x - k}{m - k_1}\right)$$

$$\approx n \frac{x - k}{m - \mu_{1k}}$$

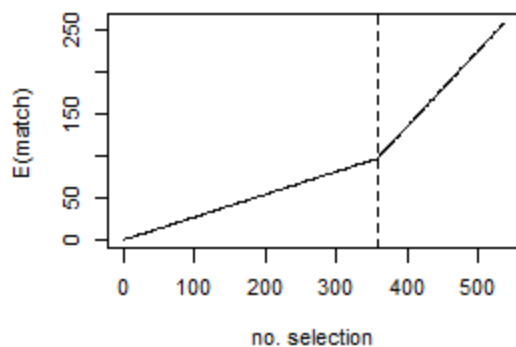
$$\mu_{1k} = E(k_1)$$

$$\dot{V}_k = (m - \mu_{1k})^{-3} n(n - m + \mu_{1k}) (x - k)(x - k - (m - \mu_{1k}))$$

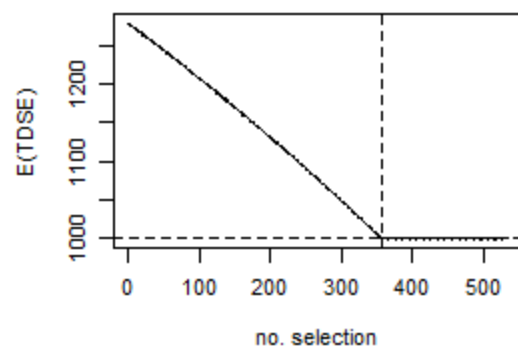
NB. had \hat{N}_k been unbiased

$$\hat{V}(\tilde{N}) = m^{-3} n(n - m) (x - r)(x - r - m)$$

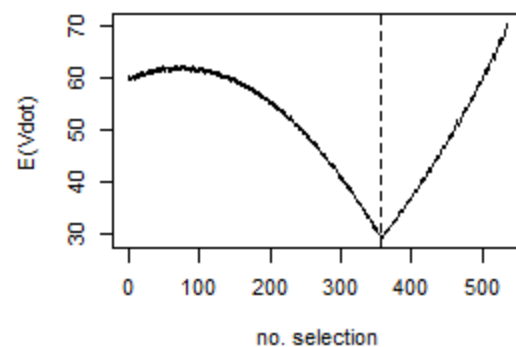
score_rate: 0.7



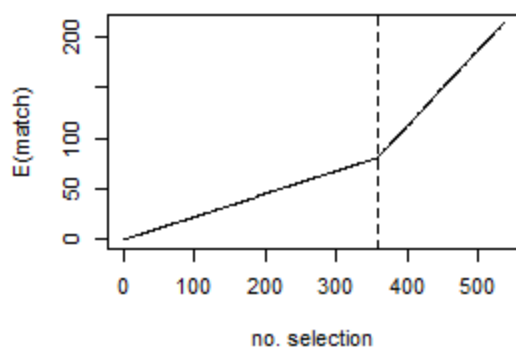
(N, n, x-r) = (1000, 900, 900)



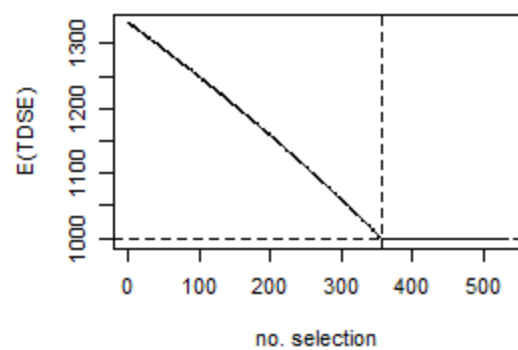
r = 250



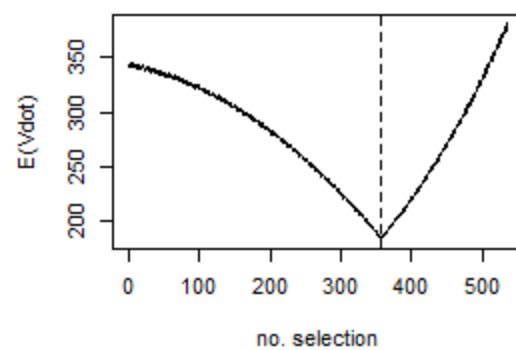
score_rate: 0.7



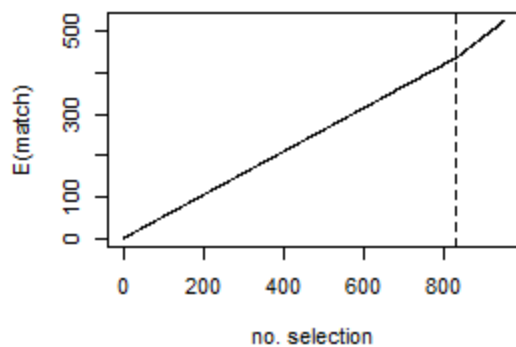
(N, n, x-r) = (1000, 750, 750)



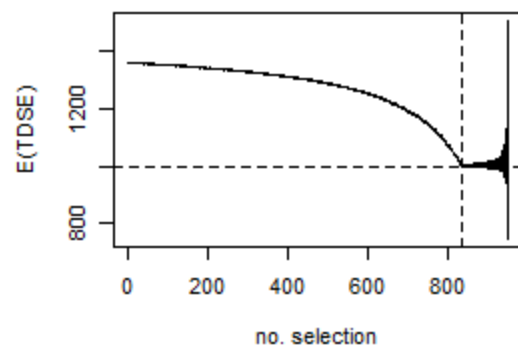
r = 250



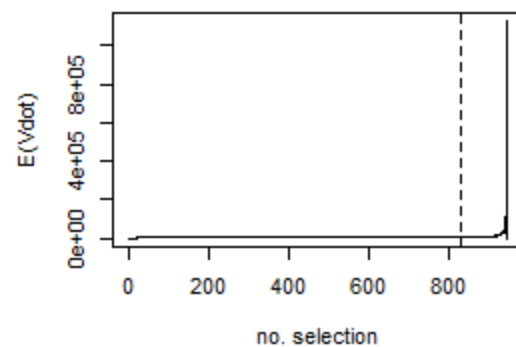
score_rate: 0.3



(N, n, x-r) = (1000, 750, 700)



r = 250



What if **both** lists have **erroneous enumerations**?

(1) Assumption of independence btw L_1 and L_2 necessary for trimmed DSE.

(2) Let the naive face-value DSE be $\hat{N}_0 = \frac{n_1 n_2}{n_{12}}$ in an obvious notation.

Let the ideal DSE be $\tilde{N} = \frac{(n_1 - r_1)(n_2 - r_2)}{n_{12} - r_{12}}$ after removing the errors.

We have $\tilde{N} < \hat{N}_0$ iff $\frac{r_{12}}{n_{12}} < \frac{r_1}{n_1} + \frac{r_2}{n_2} - \frac{r_1}{n_1} \frac{r_2}{n_2}$

i.e. provided relatively fewer erroneous records among the matched list, the face-value DSE remains biased upwards; whereas it is downwards biased

in the opposite case; and **unbiased** if $E\left(\frac{r_{12}}{n_{12}}\right) = E\left(\frac{r_1}{n_1}\right) + E\left(\frac{r_2}{n_2}\right) - E\left(\frac{r_1}{n_1}\right)E\left(\frac{r_2}{n_2}\right)$.

Notice the form $P(A)+P(B)-P(A)P(B)$ on the right-hand side, which corresponds to the case where erroneous records occur randomly across L_1 , L_2 and L_{12} ,

(3) Suppose scoring yields (k_1, k_2, k_{12}) units in L_1, L_2 and L_{12} respectively.

Put the trimmed DSE (TDSE) $\hat{N}_{k_1 k_2} = \frac{(n_1 - k_1)(n_2 - k_2)}{n_{12} - k_{12}}$

We have $\hat{N}_{k_1 k_2} < \hat{N}_0$ iff $\frac{k_{12}}{n_{12}} < \frac{k_1}{n_1} + \frac{k_2}{n_2} - \frac{k_1 k_2}{n_1 n_2}$

i.e. TDSE reduces the bias of \hat{N}_0 provide scoring is more effective than random selection, and $\tilde{N} < \hat{N}_0$ to start with.

(4) Provided $\tilde{N} > \hat{N}_0$ to start with, which is unlikely in the census context, TDSE again reduces the bias of \hat{N}_0 provide scoring is more effective than

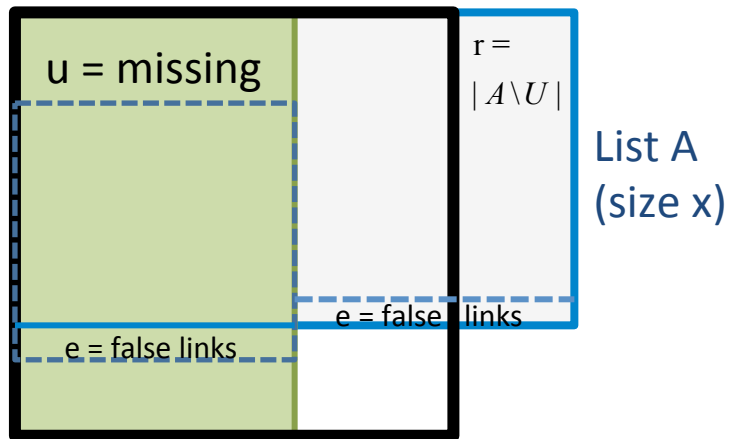
random selection. But effective means $\frac{k_{12}}{n_{12}} > \frac{k_1}{n_1} + \frac{k_2}{n_2} - \frac{k_1 k_2}{n_1 n_2}$ now.

(5) Trimming introduces bias if $\frac{r_{12}}{n_{12}} = \frac{r_1}{n_1} + \frac{r_2}{n_2} - \frac{r_1 r_2}{n_1 n_2}$ to start with, unless

scoring happens to yield $\frac{k_{12}}{n_{12}} = \frac{k_1}{n_1} + \frac{k_2}{n_2} - \frac{k_1 k_2}{n_1 n_2}$.

What about TDSE in the presence of **linkage errors**?

Survey S
(size n)



Target population U
(size N, unknown)

Black: target population U

Blue: list A with under- and over-coverage

Green: survey S with only under-coverage

Dotted area: **linked set** $S \cap A$ with linkage error

no. missing matches = u

no. false links = e (involving records in S and A)

Let the linked set be of size **M**, where $M \neq m$,

and **m** is the size of the **true matched set**.

We have $E(M|m, S, A) = m - E(u|-) + E(e|-) = m - m \cdot f + E(M|-) \cdot q$

$\Rightarrow \hat{m} = M \cdot \lambda$ where $\lambda = (1-q)/(1-f)$

$f = m^{-1}E(u|m, S, A)$ is the missing rate and $q = E(e/M|m, S, A)$ the false link rate

$\Rightarrow \tilde{N}_L = n(x-r)/(M\lambda) = \text{approx. unbiased linkage DSE (LDSE)}$

Face-value DSE: $\hat{N}_0 = nx / M$

Hypothetical ideal DSE: $\tilde{N} = n(x - r) / m$ without linkage errors

Hypothetical linkage DSE (LDSE): $\tilde{N}_L = n(x - r) / (M\lambda)$ with linkage errors

We have $\lambda = (1 - q) / (1 - f) \approx (1 - q)(1 + f)$, so that

$$\hat{N}_0 > \tilde{N}_L \Leftrightarrow r > x(1 - \lambda) \approx x[q(1 + f) - f]$$

More simply, "conservative" linkage if $q < f \Leftrightarrow \lambda > 1 \Leftrightarrow E(M|m, S, A) < m$

$$\Rightarrow \hat{N}_0 > \tilde{N}_L \approx \tilde{N}$$

It follows that given the actual linkage result, with plug-in $\hat{\lambda}$, the trimmed linkage DSE (TLDSE) can be given as

$$\hat{N}_k = n \frac{x - k}{(M - k_1)\hat{\lambda}}$$

It can be expected to behave similarly as the TDSE in the absence of linkage errors.

NB. A more subtle discussion requires one to know how λ changes under scoring.

What if trimming can not do all the job?

Need at least 2 lists, in addition to **coverage survey (S)**

		List B			
		in	out		
In U	List A	in	p_{111}	p_{110}	p_{11+}
		out	p_{101}	p_{100}	p_{10+}
			p_{1+1}	p_{1+0}	p_{1++}
		List B			
		in	out		
Out of U	List A	in	p_{011}	p_{010}	p_{01+}
		out	p_{001}	—	p_{001}
			p_{0+1}	p_{010}	

NB. subscript = domain indicator (uab); $N_{000} = 0$ by definition
 Zhang (2015) examine possible log-linear models for

- Target population U
- List-target universe U^* = union of U , A and B
- List universe U_L = union of A and B

Two well-defined, generalisable models

Model (10): put $\theta_{ab} = P(I_{iu} = 0 \mid I_{iab} = 1)$

unsaturated log-linear model for list-target universe

$$\text{logit } \theta_{11} = \text{logit } \theta_{10} + \text{logit } \theta_{01}$$

approximately the same as

$$\log \theta_{11} = \log \theta_{10} + \log \theta_{01} \quad \Leftrightarrow \quad \theta_{11} = \theta_{10} \theta_{01}$$

$$\Leftrightarrow P(i \notin U \mid i \in A \text{ and } i \in B) = P(i \notin U \mid i \in A \setminus B) \cdot P(i \notin U \mid i \in B \setminus A)$$

Model (11):

$$\log \theta_{11} = \log \theta_{1+} + \log \theta_{+1} \quad \Leftrightarrow \quad \theta_{11} = \theta_{1+} \theta_{+1}$$

$$\Leftrightarrow P(i \notin U \mid i \in A \text{ and } i \in B) = P(i \notin U \mid i \in A) \cdot P(i \notin U \mid i \in B)$$

Zhang (2015) refers this as **pseudo conditional independence (PCI)**

NB. for any random events Z , X and Y

$$\text{CI:} \quad P(X \cap Y \mid Z) = P(X \mid Z) \cdot P(Y \mid Z)$$

$$\text{PCI:} \quad P(Z \mid X \cap Y) = P(Z \mid X) \cdot P(Z \mid Y)$$

Generalisation:

Generic log-linear models for error- or hit-rates of 3-list capture table

Model Restriction	Model Interpretation
-	Saturated model
$\alpha_{123} = 0$	Null 2nd-order PCI-interaction among (L_1, L_2, L_3)
$\alpha_{12} = 0$	PCI between L_1 and L_2
$\alpha_{123} = -\alpha_{12}$	Conditional PCI between L_1 and L_2 given L_3
$\alpha_{12} = \alpha_{123} = 0$	Both PCI and conditional PCI between L_1 and L_2
$\alpha_{12} = \alpha_{13} = \alpha_{123} = 0$	PCI between L_1 and (L_2, L_3)
$\alpha_{12} = \alpha_{13} = \alpha_{23} = \alpha_{123} = 0$	Mutual PCI between L_1, L_2 and L_3

Back to the 2-list situation...

Table 3: Range of relative bias under model (10) and (11) for census enumeration error adjustment. Census enumeration = 1000, register enumeration = 1200, census-register enumeration = 900. Error rate of census θ_{1+} , register enumeration (θ_{+1}) , census-register enumeration (θ_{11}) , where $0 < \theta_{11} < \theta_{1+}$. All numbers in %.

Model (10)	Register error rate			
Census error rate	1	5	10	20
0.2	(0.078, 0.078)	(-0.11, -0.11)	(-0.48, -0.48)	(-3.4, -3.4)
0.5	(-0.038, 0.43)	(-0.88, 0.32)	(-2.5, 0.095)	(-16, -1.6)
1	(-0.25, 1)	(-2.3, 1)	(-6.3, 1)	(-38, 1)

Model (11)	Register error rate			
Census error rate	1	5	10	20
0.2	(0.11, 0.11)	(0.11, 0.11)	(0.1, 0.1)	(0.089, 0.089)
0.5	(0.11, 0.45)	(0.091, 0.44)	(0.068, 0.44)	(0.014, 0.43)
1	(0.1, 1)	(0.065, 1)	(0.012, 1)	(-0.11, 1)

Combine Trimming e.g. with Model (11):

Trim the initial lists A and B in order to move them into the well-fitting range of model (11) ?

The Irish case: data sources

Sources (Cradle to Grave)

- Births
- Child benefit
- Post primary
- Higher Ed
- Further Ed
- Employee
- Self Employed
- Social Welfare
- Drivers
- Medical (to be)
- Pensions
- Deaths



PAR - Person Activity Register
(CSO-PPSN)
2011 (4.3m)
Admin Active
Population during
year
Signs of life (SoL)

Person Activity Register (PAR) only includes those for whom there has been evidence of engaging in **key** administrative systems in a given year.

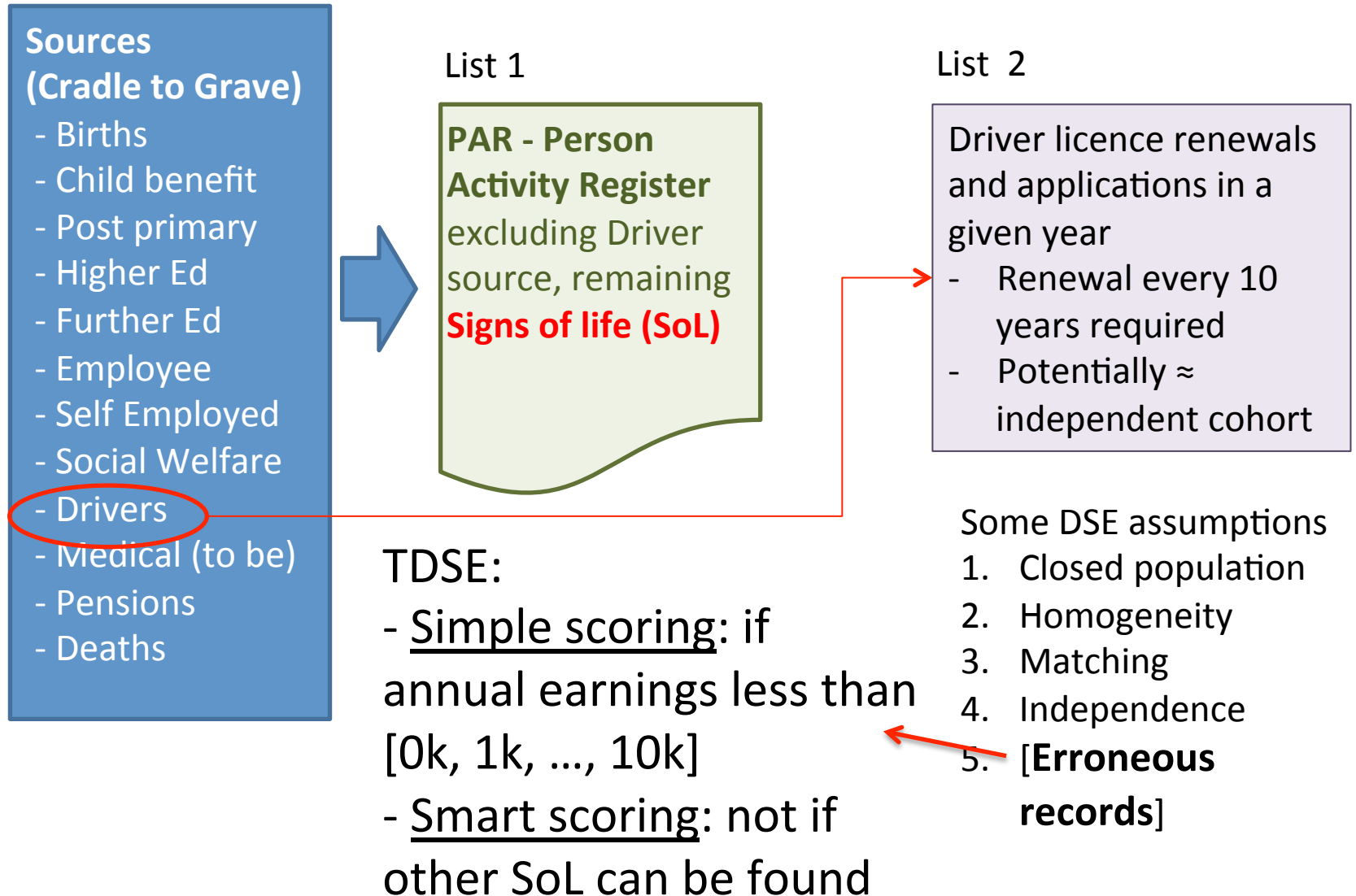
- **Over-coverage** can be “**defined away**”, but is subjected to discussion.
- **Under-coverage** by and large **to be expected**

Population estimate = PAR enumeration?

For the year **2011** in Ireland,
taking a “**signs of life**” approach
the Person Activity Register identified activity from
4.35m persons
on Public Administration Systems.

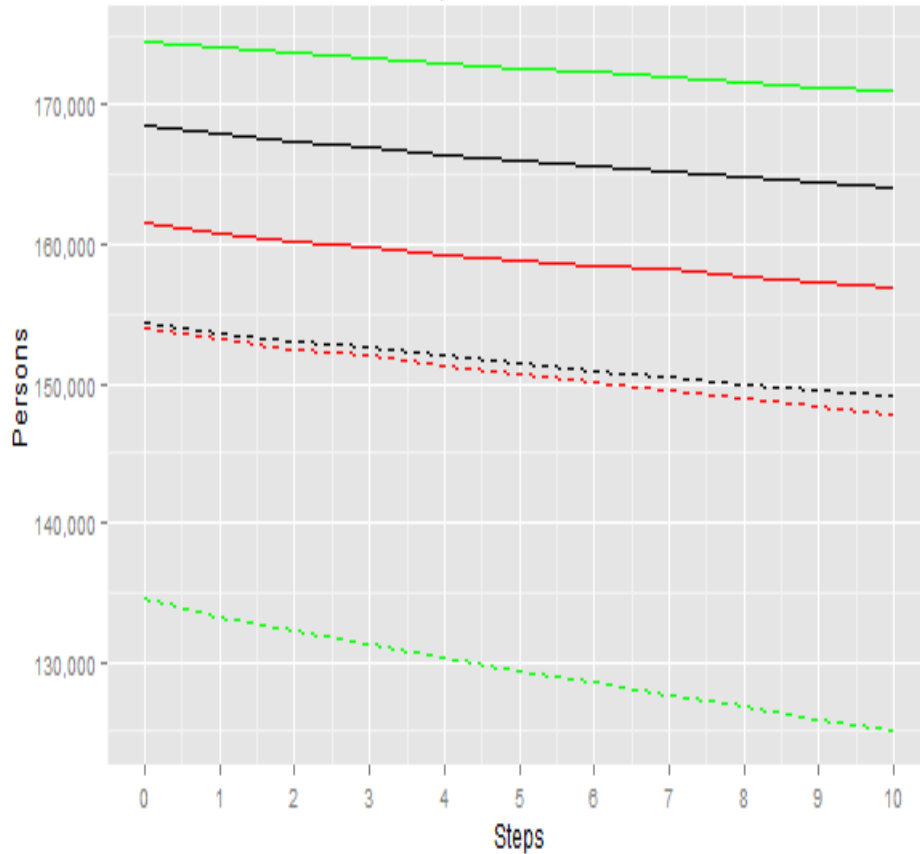
4.59m people were enumerated
on Census night in 2011

The Irish case: create a setting for DSE?

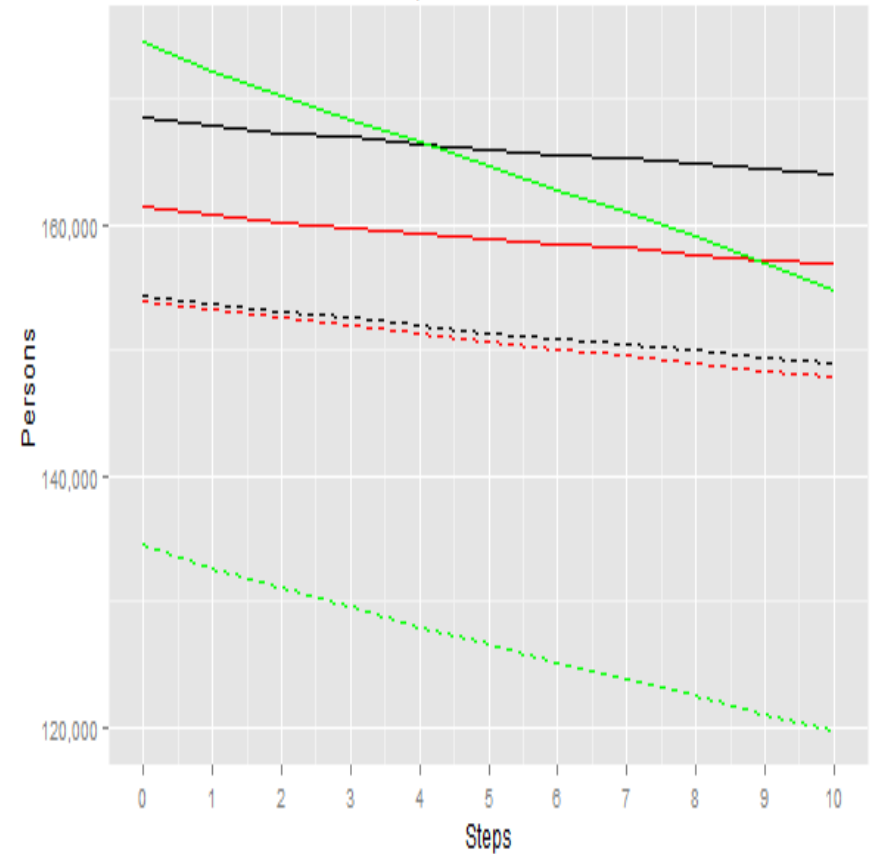


Age 40 - 44 PAR **scoring**: smart (left), simple (right)

Directly observed population: Males (dashed) and Female (continuous), 40 to 44
.... Red 2009, Black 2011 and Green 2013



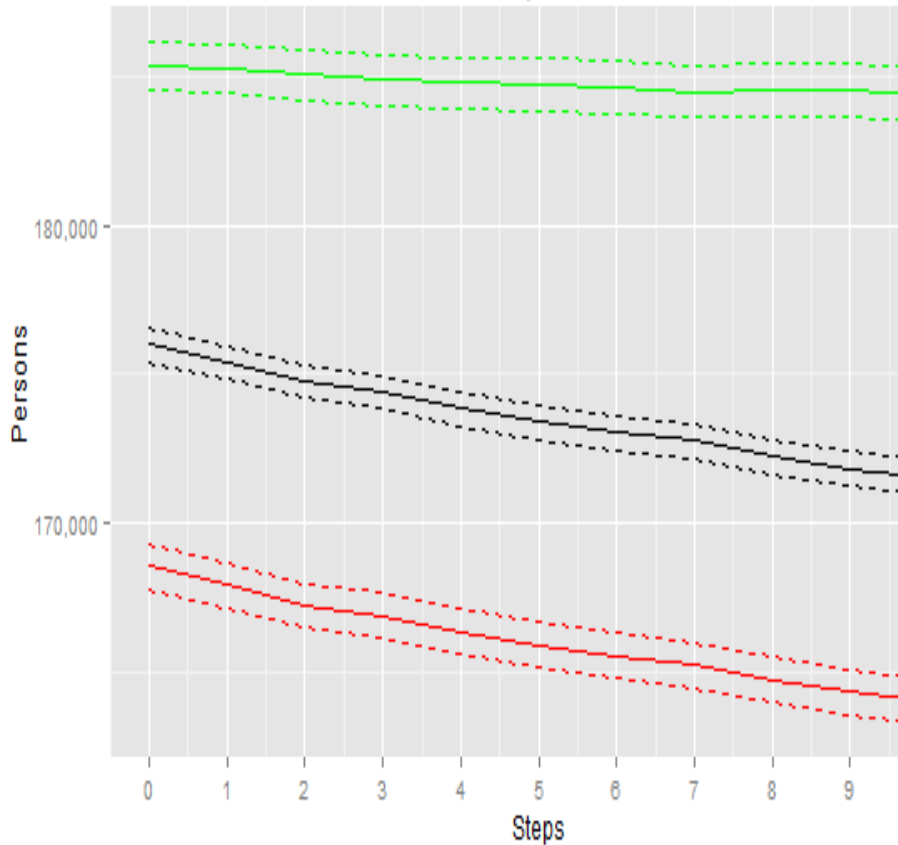
Directly observed population: Males (dashed) and Female (continuous), 40 to 44
.... Red 2009, Black 2011 and Green 2013



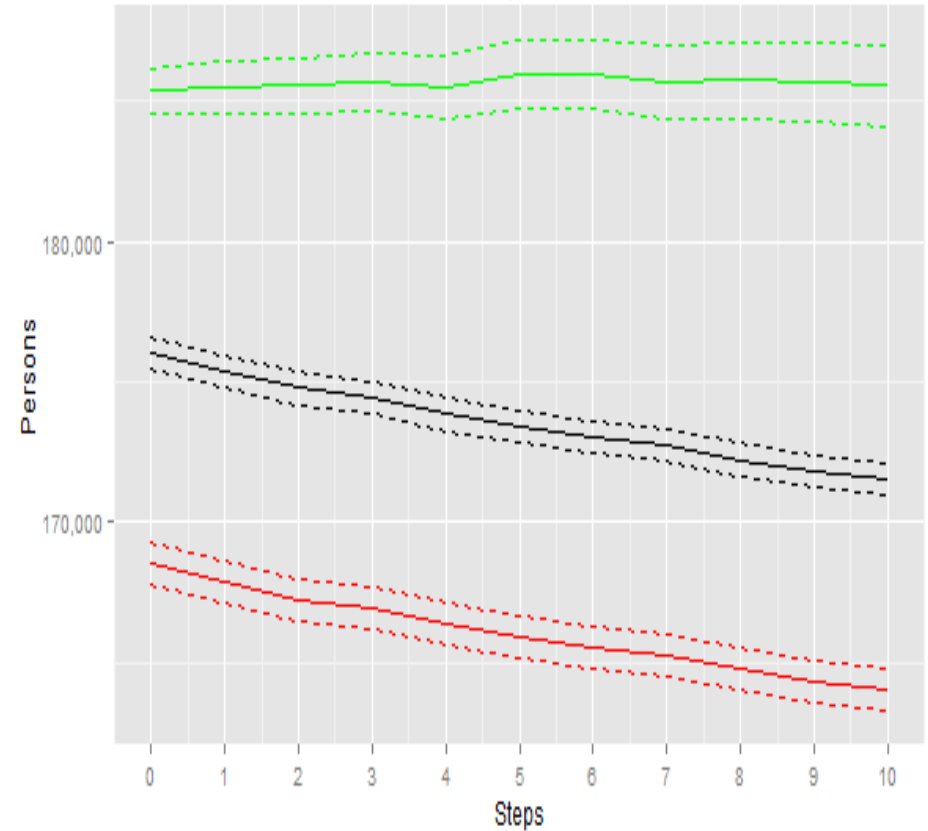
2013: PAR not finalised, lack of self-employed

Female 40 - 44 PAR **TDSE**: smart (left), simple (right)

Population estimates: Females, 40 to 44
.... with 95% CIs Red 2009, Black 2011 and Green 2013



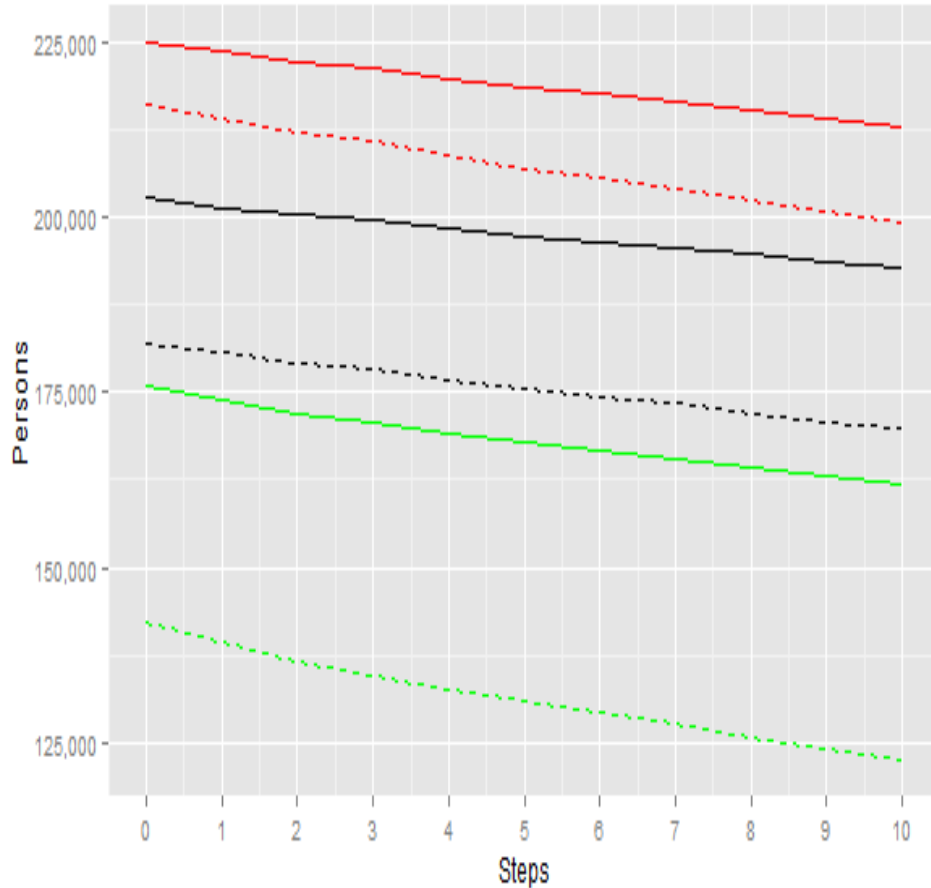
Population estimates: Females, 40 to 44
.... with 95% CIs Red 2009, Black 2011 and Green 2013



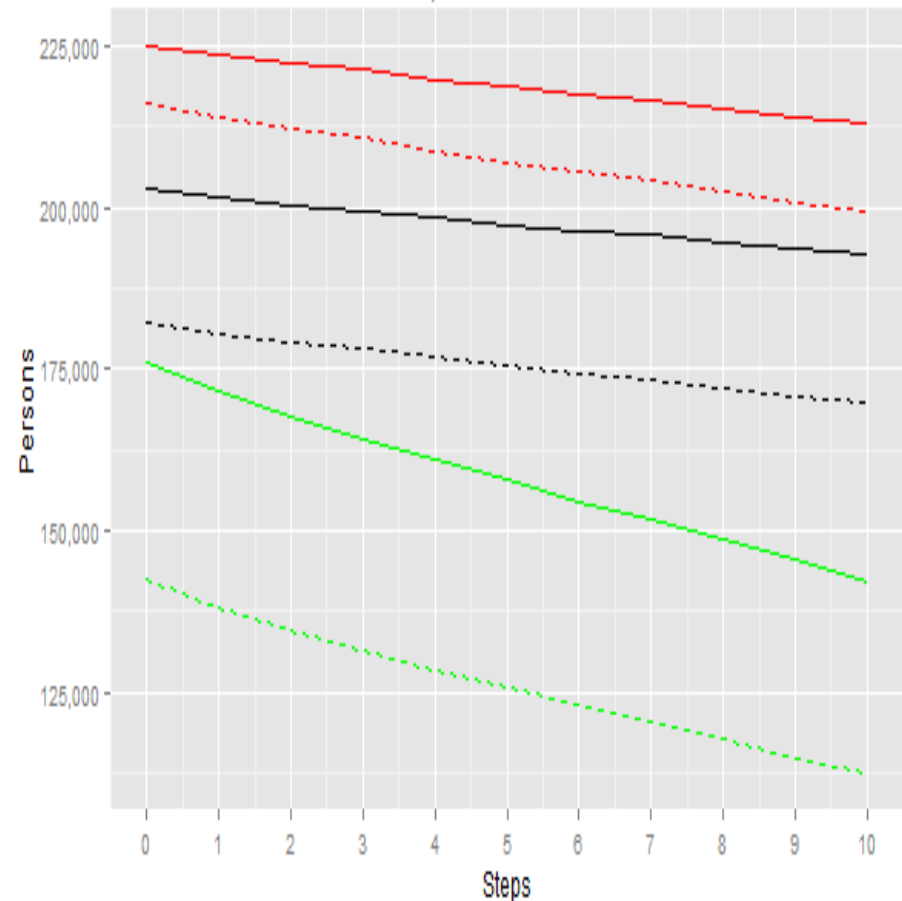
2013: PAR not finalised, lack of self-employed

Age 25 - 29 PAR **scoring**: smart (left), simple (right)

Directly observed population: Males (dashed) and Female (continuous), 25 to 29
.... Red 2009, Black 2011 and Green 2013



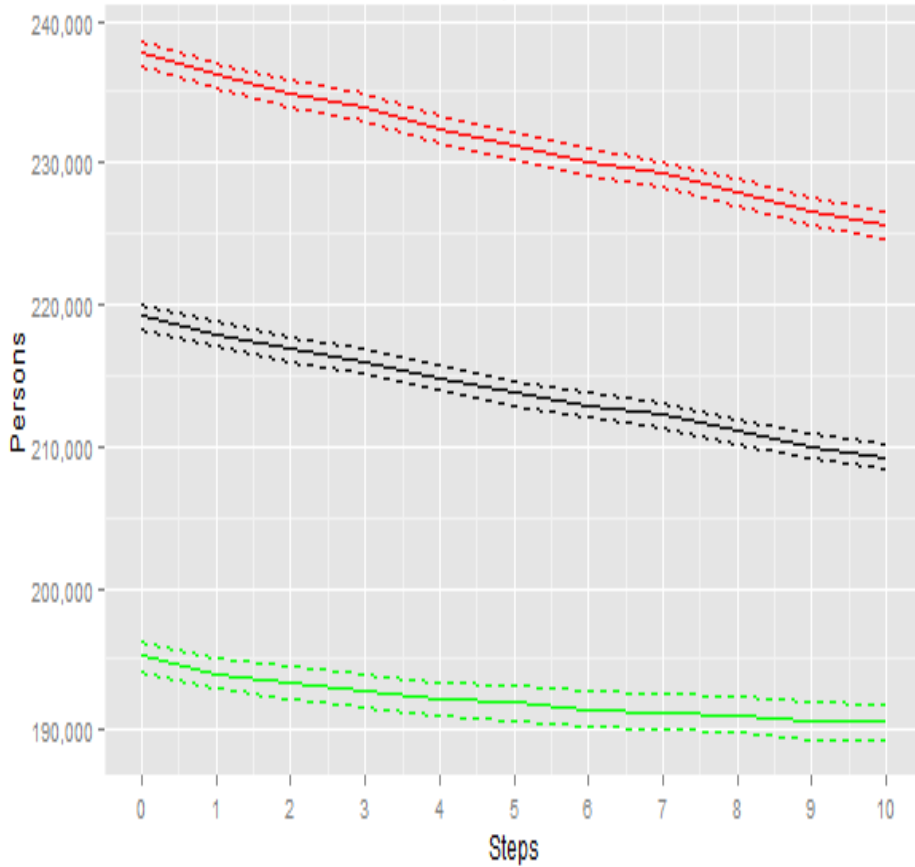
Directly observed population: Males (dashed) and Female (continuous), 25 to 29
.... Red 2009, Black 2011 and Green 2013



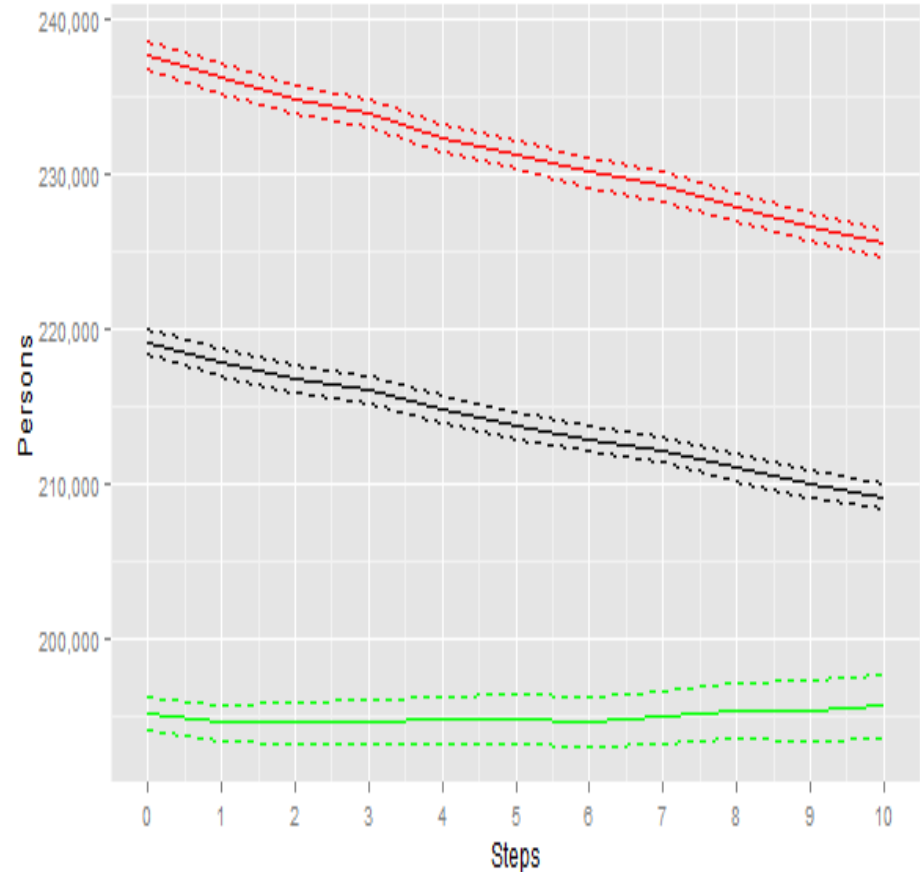
2013: PAR not finalised, lack of self-employed

Female 25 - 29 PAR **TDSE**: smart (left), simple (right)

Pop estimates: Females, 25 to 29
.... with 95% CIs Red 2009, Black 2011 and Green 2013



Pop estimates: Females, 25 to 29
.... with 95% CIs Red 2009, Black 2011 and Green 2013



2013: PAR not finalised, lack of self-employed

Discussion

- As by theory: larger $k \neq$ smaller N estimate
 - Bias: Scoring rate vs. prevalence of err. records
 - Variance: increases with extent of scoring
- Stability & composition of “Driver cohorts”:
 - had everyone become a driver only at exactly the same point in their life, only a sub-population can be proved in any given year
- Implications of $L_2 =$ Driver cohorts
 - Target population = one for which Driver cohorts have only under- but not over-coverage
 - Temporary workers or other special groups: random but not systematic under-enumeration
- Address/Immobility Register and address registration:
Critical for census-like detailed population statistics