

The 'history' of the R-indicator

Jelke Bethlehem
Fannie Cobben
Barry Schouten



First RISQ-meeting, 7-8 April 200

In this presentation...

- Statistics Netherlands
- Research programme
- R-indicators; theoretical background
- R-indicators in practice
- Discussion and future research



Primary data collection

- Data collection modes:
 - Face-to-face (CAPI)
 - Telephone (CATI)
 - Web
 - Paper
- All surveys based on probability samples from municipality registers
- Registered land-line phone numbers are linked from commercial databases (70% coverage)
- Web data collection only in pilot studies using letters + logins to secured website (80% coverage)
- At present no household survey employs a mixed-mode design



Secondary data collection

Statistics Netherlands Act: By law 'allowed' to use government registers and administrative data as input to the production of statistics

Examples:

- Municipality registers (Population register)
- Tax Board registers on wages, VAT, profits, incomes
- Registers for various government allowances
- Register on value of real estate

Population register functions as backbone to both probability samples and other government registers



Strategic Programme *Nonresponse, Difficult Groups and Mixed-mode*

Research projects:

1. Nonresponse reduction
2. Nonresponse adjustment
3. Difficult groups
4. Mixed-mode data collection



Response enhancement

- Differentiated data collection protocols
- Responsive/adaptive designs

Indicators for representative response (R-indicators)

Indicators as tools to:

- compare surveys in time
- compare different data collection strategies
- monitor and control data collection

Consequence: Focus on response behavior, i.e. independent of survey items.

Important: Auxiliary information and paradata are crucial to any indicator. An indicator must always be published together with the available external information.



Representativity; what?

Stoop (2005):

- There is no such thing as a representative sample

Schnell (1997):

- '*Representative sampling*' is an immeasurable, non-scientific concept, without any specific meaning

Kruskal en Mosteller (1979):

- 9 definitions of representativity
- Recommendation: do not use the word 'representative', but specify what you mean by it



R-indicators: Definition and Concept

Definition (strong): *A response subset is representative with respect to the sample if the response propensities are the same for all units in the population and if the response of a unit is independent of the response of all other units.*

Definition (weak): *A response subset is representative for a categorical variable X if the average response propensity over the categories of X is constant.*



Notation

Response probabilities:

Population

$$N, i = 1, 2, \dots, N$$

Selection - indicator

$$s_i = \begin{cases} 1, & \text{if person } i \text{ selected in sample} \\ 0, & \text{else} \end{cases}$$

Sample size

$$n = \sum_{i=1}^N s_i$$

First order inclusion probabilities $\pi_i = P(s_i = 1)$

Response - indicator

$$r_i = \begin{cases} 1, & \text{response person } i \\ 0, & \text{nonresponse person } i \end{cases}$$

Response probability

$$\rho_i = P(r_i = 1 | s_i = 1) = \frac{\exp(-\beta_i X_i)}{1 + \exp(-\beta_i X_i)}$$



R-indicators – Example

Variation of response propensities in population

$$R(\tilde{\rho}) = 1 - 2 \sqrt{\frac{1}{N-1} \sum_{i=1}^N (\rho_i - \bar{\rho})^2}$$

Estimated variation of response propensities

$$\hat{R}(\tilde{\rho}) = 1 - 2 \sqrt{\frac{1}{N-1} \sum_{i=1}^N \frac{s_i}{\pi_i} (\rho_i - \bar{\rho}_{HT})^2}$$

Estimated variation of estimated response propensities

$$\hat{R}(\hat{\tilde{\rho}}) = 1 - 2 \sqrt{\frac{1}{N-1} \sum_{i=1}^N \frac{s_i}{\pi_i} (\hat{\rho}_i - \hat{\rho}_{HT})^2}$$

R-indicators – Features

Interpretation: Dependence on X 's and n

Normalization of R-indicators: Relate to non-response bias and RMSE under worst case scenario

$$|B(\hat{y}_{HT})| \leq \frac{S(\tilde{\rho})S(y)}{\hat{\rho}} \leq \frac{S(y)(1 - R(\tilde{\rho}))}{2\hat{\rho}}$$

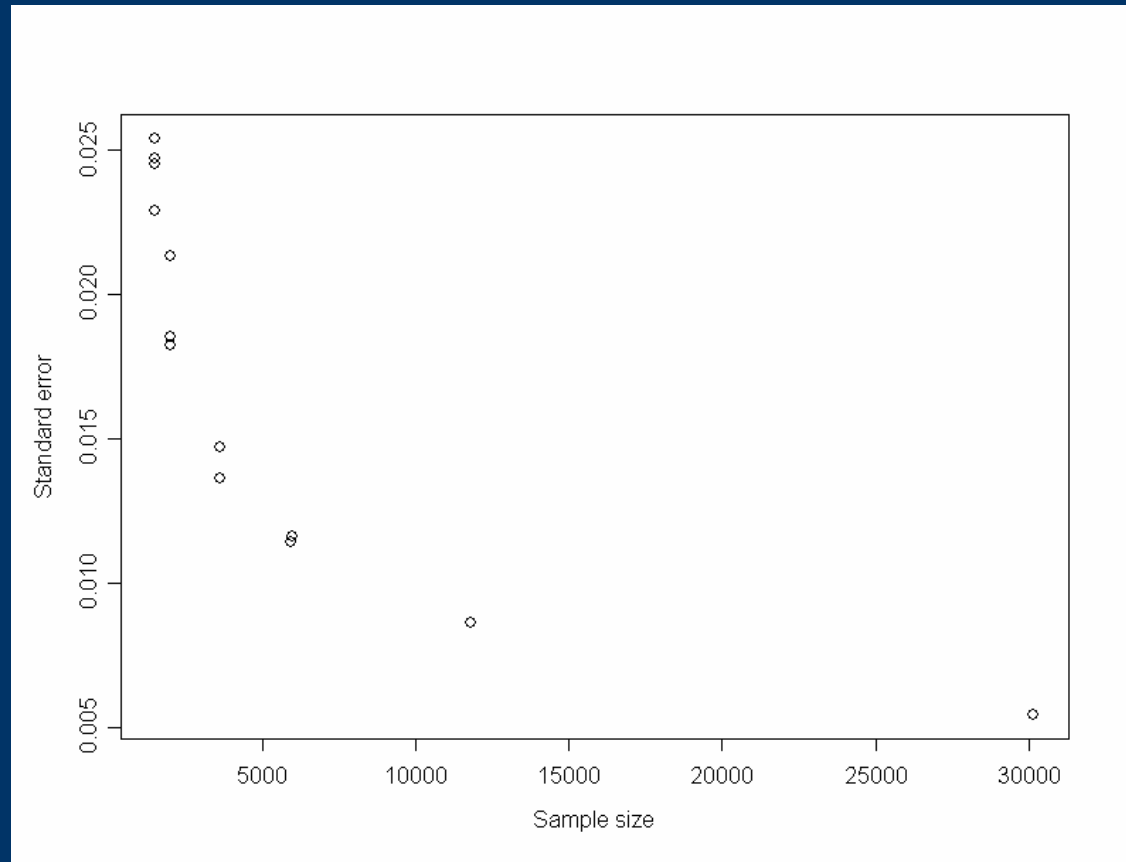
$$R(\tilde{\rho}) \geq 1 - 2 \frac{\hat{\rho}\gamma}{S(y)}$$

$$\begin{aligned} RMSE(\hat{y}_{HT}) &= \sqrt{B^2(\hat{y}_{HT}) + Var(\hat{y}_{HT})} \\ &\leq \sqrt{B^2(\hat{y}_{HT}) + \left(1 - \frac{n\bar{\rho}}{N}\right) \frac{S^2(y)}{n\bar{\rho}}} \end{aligned}$$

$$R(\tilde{\rho}) \geq 1 - \frac{2\hat{\rho}}{S(y)} \sqrt{\gamma^2 - \left(1 - \frac{n\hat{\rho}}{N}\right) \frac{1}{4n\hat{\rho}}}$$



R-indicators – Features

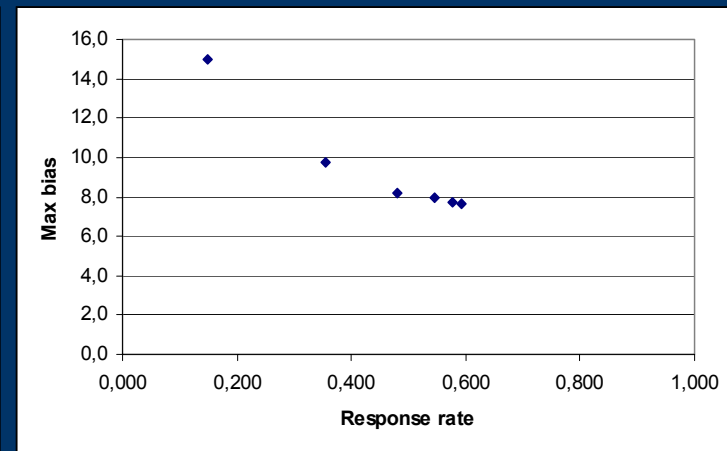
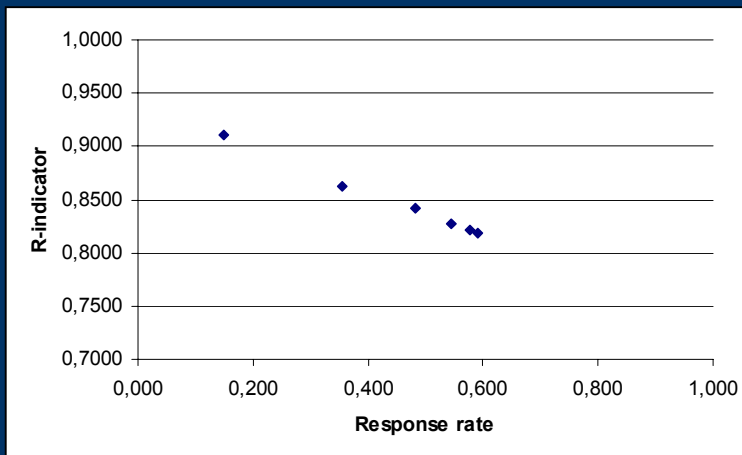


Example – Contact Attempts

Survey POLS 1998, sample size $n = 35.893$

CAPI in first month, CATI in second month

$X = \text{Age, ethnic group, region}$



Example: Call Back & Basic Question

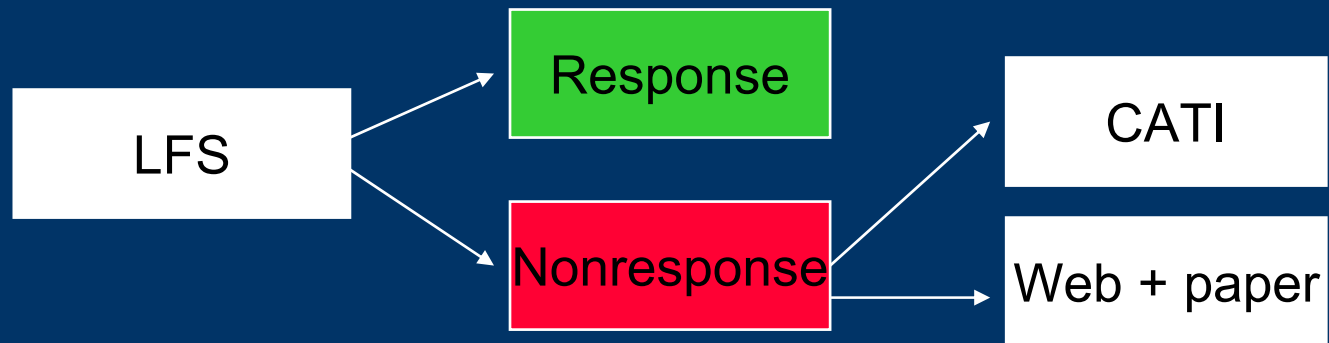
Survey LFS July – October 2005

Call-back approach (Hurwitz 1949)

- Selection of best performing interviewers
- Additional training of interviewers
- Incentives
- Paper summaries of household characteristics

Basic-question approach (Kersten & Bethlehem 1988)

Condensed questionnaires in CATI, paper, web



Example: Call Back & Basic Question

LFS n=18.076, CBA n=785

X=phone, region, ethnic group, household type, urbanity

	Response	R-indicator	Max bias
LFS	62,2%	80,1%	8,0%
LFS + CBA	76,9%	85,1%	4,8%

LFS n=18.076, BQA n=942

X=household type, urbanity, age, gender, job, allowance

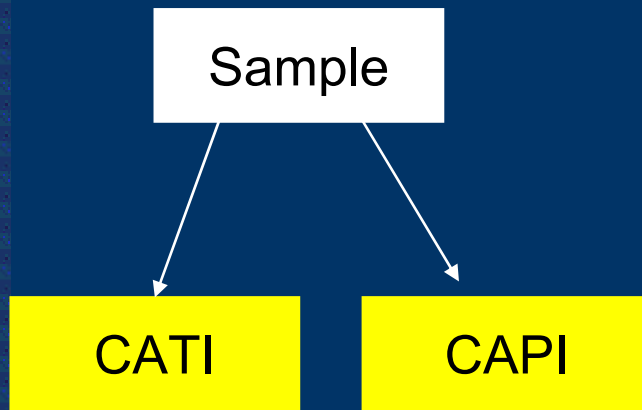
	Response	R-indicator	Max bias
LFS	62,2%	80,1%	8,0%
LFS, phone	68,5%	86,3%	5,1%
LFS + CBA	75,6%	78,0%	7,3%
LFS + CBA, phone	83,0%	87,5%	3,8%



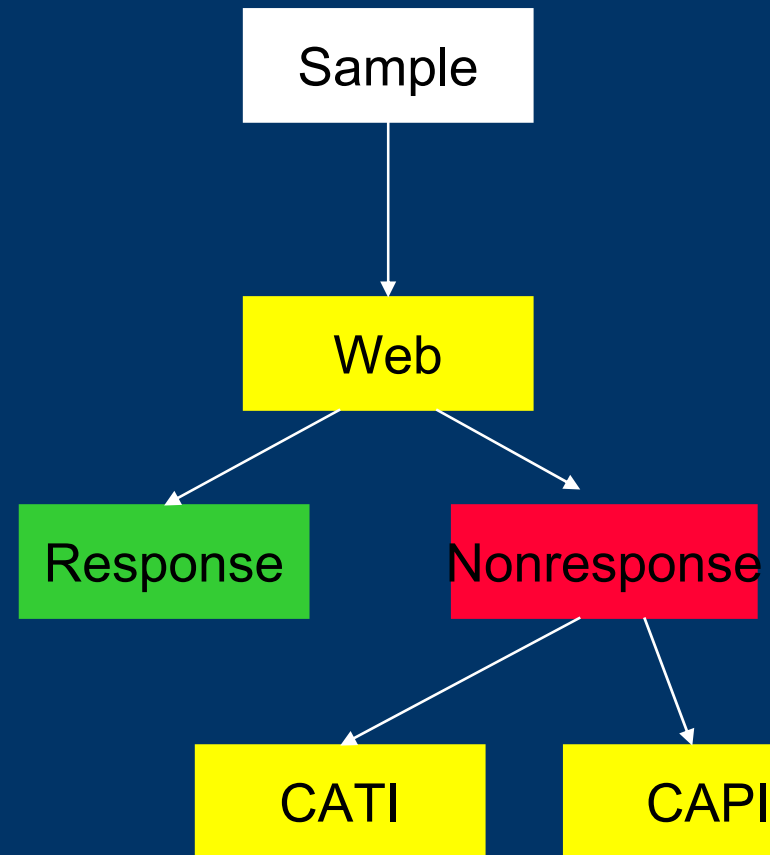
Example: Mixing Modes (1)

Safety Monitor 2006

Reference survey



Pilot survey



Example: Mixing Modes (1)

Safety Monitor 2006

X=urbanity, household type, ethnic group, age

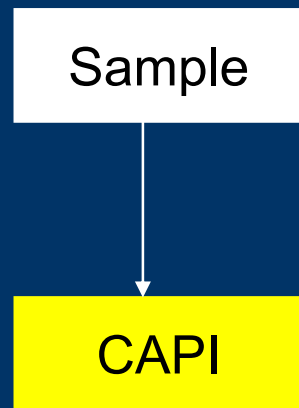
	n	Response	R-indicator	Max bias
Reference	30.139	68,9%	81,4%	6,8%
Pilot, web	3.615	30,2%	77,8%	18,4%
Pilot, total	3.615	64,7%	81,2%	7,3%



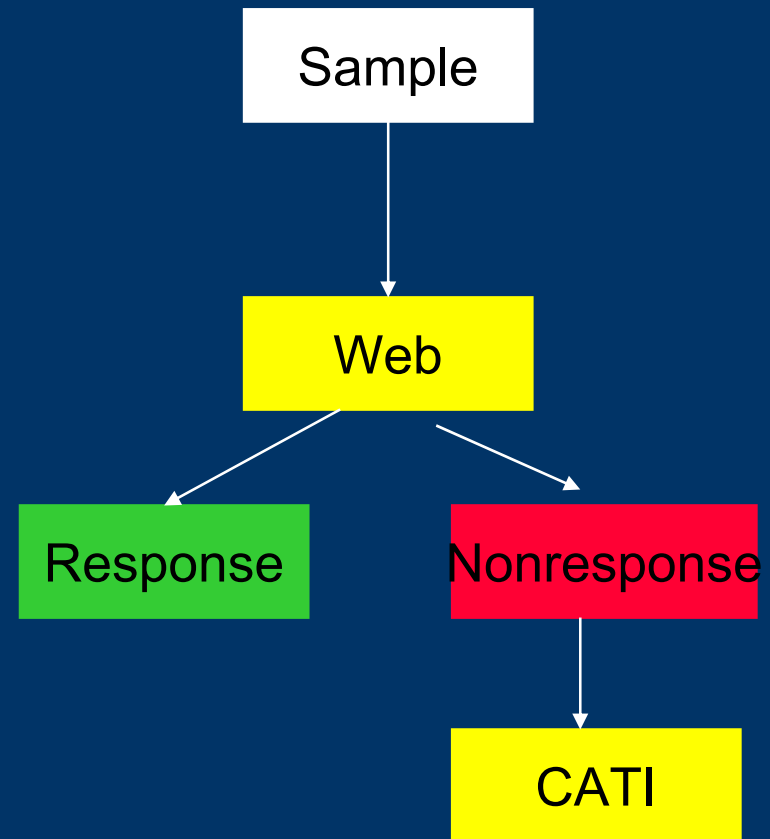
Example: Mixing Modes (2)

Informal Economy 2006

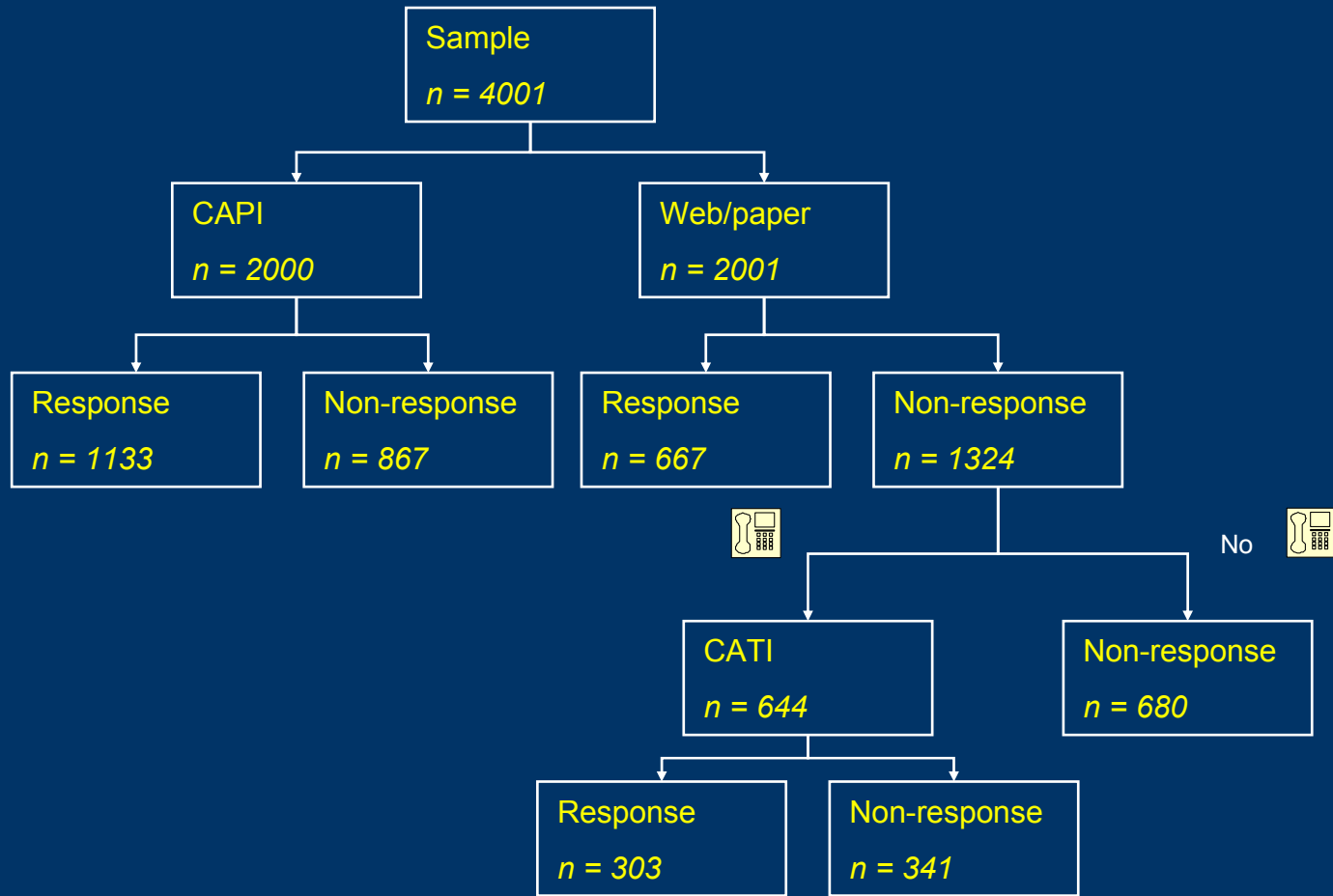
Pilot survey 1



Pilot survey 2



Pilot Informal Economy 2006



Example: Mixing Modes (2)

Informal Economy 2006

X= urbanity, household type, ethnic group, age

	n	Response	R-indicator	Max bias
CAPI	2.000	56,7%	77,2%	10,1%
Web	2.001	33,8%	85,1%	11,0%
Web + CATI	2.001	49,0%	78,0%	11,2%



Example: Incentives

Survey LFS 2005

Incentives:

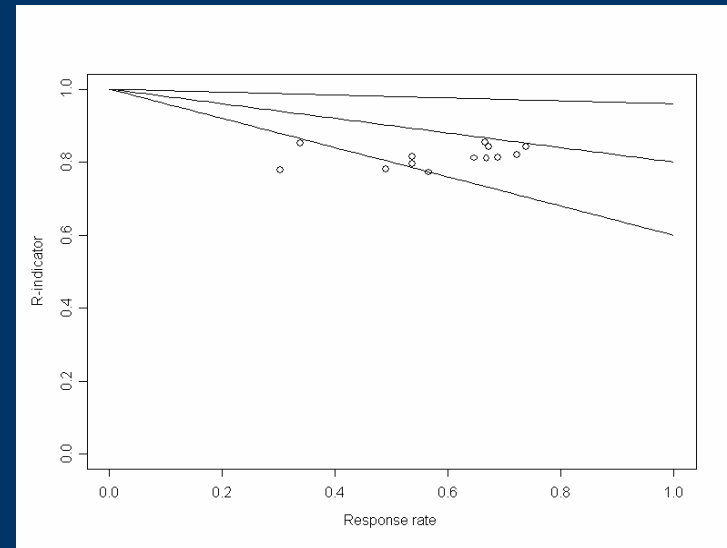
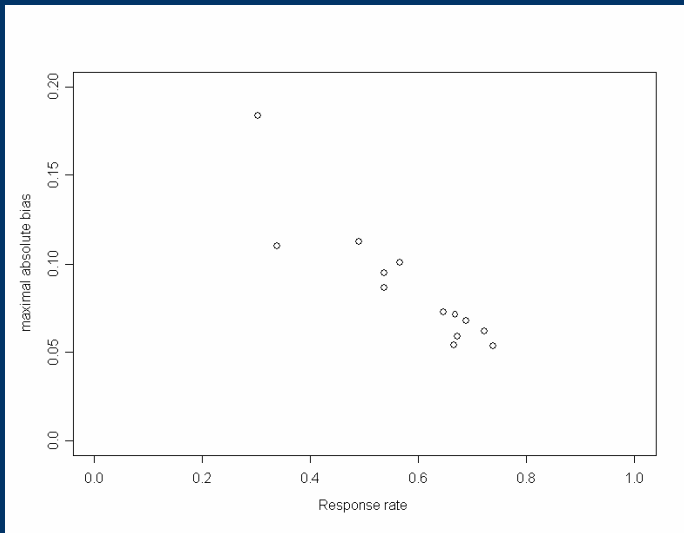
1) no stamps, 2) 5 stamps , and 3) 10 stamps

X= urbanity, average house value, ethnic group, size of household

	n	Response	R-indicator	Max bias
No	11.774	66,6%	85,5%	5,4%
5	5.906	72,2%	82,1%	6,2%
10	5.982	73,8%	84,2%	5,4%



Example: Maximal bias



Discussion & future research

- Can we ignore survey items?
- Are there alternative R-indicators?
- Can R-indicators be tools in monitoring or even controlling survey data collection?
- Can R-indicators help in comparing different surveys (possibly over time)?
- How to interpret the values of R-indicators?



Discussion & future research

Short term:

- Extend theory to situation where only population
- totals are available
- Construction of R-indicator confidence intervals

Longer term:

- RISQ
- Responsive designs

