Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Non-probability sampling and big data

EMOS Webinar 2022/23

Ralf Münnich

18. January 2023

EUROPEAN
MASTER IN
OFFICIAL
STATISTICS
EMOS

UNIVERSITÄT
**TRIER**

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# What is the problem with survey data?

▶ Gallup opinion poll in 1948 on US elections
  ▶ Dewey (Republicans) versus Truman (Democrats)
  ▶ Use of quota sampling
  ▶ Prediction: Dewey – but survey stopped early
  ▶ Winner: Truman

▶ Johnson's red bus (Brexit), Trump election, etc.

▶ Huge debate in Germany:
  Market and opinion research versus internet surveys
  Non-response versus web selectivity

▶ What is a really good survey?

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# What is the problem with survey data?

▶ Gallup opinion poll in 1948 on US elections
  ▶ Dewey (Republicans) versus Truman (Democrats)
  ▶ Use of quota sampling
  ▶ Prediction: Dewey – but survey stopped early
  ▶ Winner: Truman

▶ Johnson's red bus (Brexit), Trump election, etc.

▶ Huge debate in Germany:
  Market and opinion research versus internet surveys
  Non-response versus web selectivity

▶ What is a really good survey?

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# What is the problem with survey data?

- ▶ Gallup opinion poll in 1948 on US elections
  - ▶ Dewey (Republicans) versus Truman (Democrats)
  - ▶ Use of quota sampling
  - ▶ Prediction: Dewey – but survey stopped early
  - ▶ Winner: Truman

- ▶ Johnson's red bus (Brexit), Trump election, etc.

- ▶ Huge debate in Germany:
  Market and opinion research versus internet surveys
  Non-response versus web selectivity

- ▶ What is a really good survey?
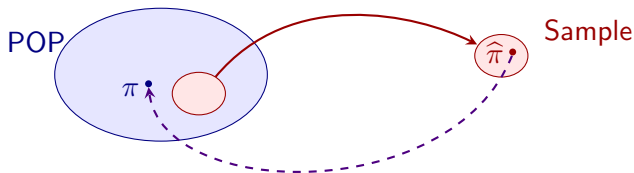
Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Content

Statistical inference and quality

Compensation methods

Web surveys and big data

Conclusion and outlook

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik
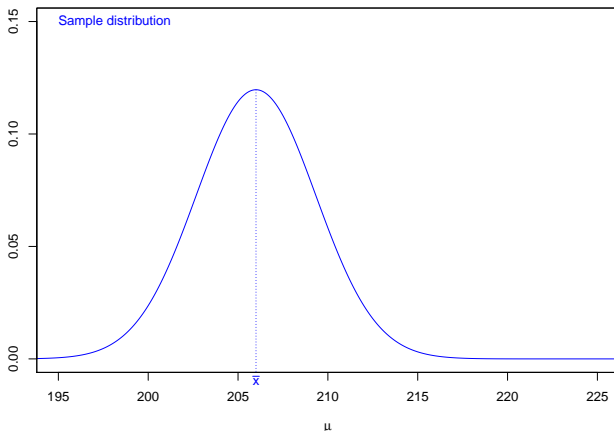
## General idea of *estimation*

We are interested in population parameters which are generally unknown (here: $\pi$).

After analysing populations using methods of descriptive statistics, we now draw a sample of the population and evaluate the outcome ($\rightarrow$ point estimation).
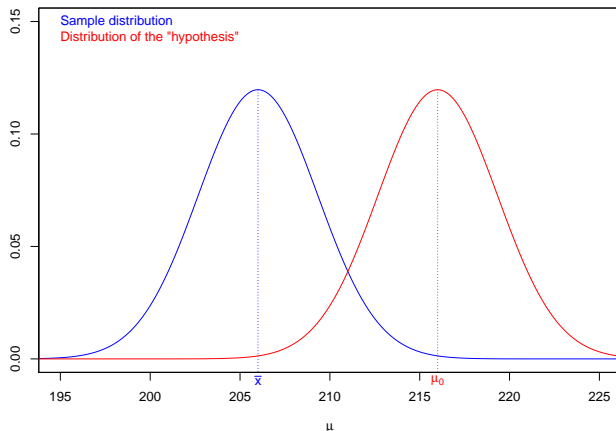


Additionally, we want to specify an interval of *plausible* values ($\rightarrow$ interval estimation in terms of providing quality information).

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# General idea of inference



▶ Analysis of sample (estimation, e.g. using $\overline{x}$)
▶ Hypothesis for population (e.g. $\mu_0$ for $\mu$)
▶ True distribution *unknown* in reality

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik
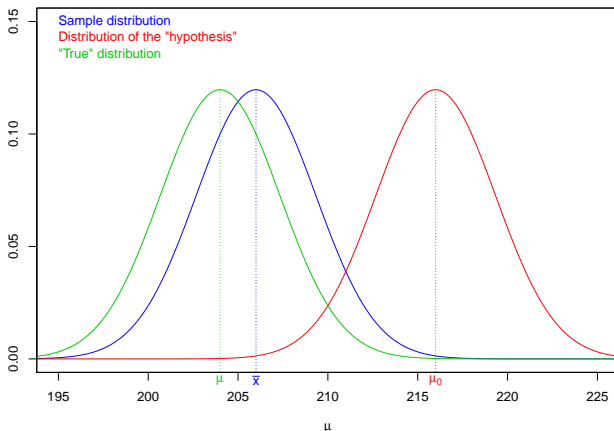
# General idea of inference



- ▶ Analysis of sample (estimation, e.g. using $\overline{x}$)
- ▶ Hypothesis for population (e.g. $\mu_0$ for $\mu$)
- ▶ True distribution *unknown* in reality

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# General idea of inference



- ▶ Analysis of sample (estimation, e.g. using $\overline{x}$)
- ▶ Hypothesis for population (e.g. $\mu_0$ for $\mu$)
- ▶ True distribution *unknown* in reality

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# What is the impact of a quality concept?

Relevance of the statistical concept:
   End-user, *user needs*, hierarchical structure and contents

Accuracy and reliability:
   ▶ Sampling errors: standard error, CI coverage
   ▶ Non-sampling errors: nonresponse, coverage error, measurement errors

Timeliness and punctuality: Time and duration from data acquisition until publication

Coherence and comparability: Preliminary and final statistics, annual and intermediate statistics (regions, domains, time)

Accessibility and clarity: Publication of data, analysis and method reports

Completeness

See: European Statistics Code of Practice

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## How does this fit in a data science context?

Are all requirements of the previous slides met in a general framework of data collection in data science?

▶ Is the frame (core population) known and well-addressed?

▶ How is the data gathering process controlled?
   Is every unit separately drawable in a completely known way?

▶ Does every unit provide full information?
   Non-response is not solely an issue in public surveys!

▶ Are there any other sources of imprecision?
   Is the measuring process adequate/precise?

Note: possible sources of imprecision have to be controlled, and if possible, measured. The output has to be evaluated in light of the **data gathering process** including all these drawbacks. This includes also possible *corrections* of the results.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Non–probability samples

**Main differences to classical probability samples:**

▶ Uncontrolled / non–random data–generating process

▶ Missing 'sampling'–information
Inclusion / participation probability $\pi_i$ unknown
(and also $\pi_{ij}$)

▶ Coverage of the target population not assured
$\pi_i = 0$ possible      (possibly also overcoverage)

▶ Possibly poor representativity & selection bias

▶ But: usually lower costs (probably also faster)
Easier to obtain certain variables

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

| Data source | Design weights $w$ | Calibration weights $d$ | Calibration variables $X$ | | | Target variables $Y$ | | | Response variables $Z$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Non–probability sample ($n$) | ? | ? | $x_{11}^n$ | $\cdots$ | $x_{1p}^n$ | $y_{11}^n$ | $\cdots$ | $y_{1p}^n$ | $z_{11}^n$ | $\cdots$ | $z_{1p}^n$ |
| | | | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| | | | $x_{n^n1}^n$ | $\cdots$ | $x_{n^np}^n$ | $y_{n^n1}^n$ | $\cdots$ | $y_{n^np}^n$ | $z_{n^n1}^n$ | $\cdots$ | $z_{n^np}^n$ |
| Calibration target data ($c$) | $w_1^c$ $\vdots$ $w_{n^c}^c$ | $d_1^c$ $\vdots$ $d_{n^r}^c$ | $x_{11}^c$ $\vdots$ $x_{n^c1}^c$ | $\cdots$ $\ddots$ $\cdots$ | $x_{1p}^c$ $\vdots$ $x_{n^cp}^c$ | ? | | | | | |
| Response–reference data ($r$) | $w_1^r$ $\vdots$ $w_{n^r}^r$ | $d_1^r$ $\vdots$ $d_{n^r}^r$ | | | | ? | | | $z_{11}^r$ $\vdots$ $z_{n^r1}^r$ | $\cdots$ $\ddots$ $\cdots$ | $z_{1p}^r$ $\vdots$ $z_{n^rp}^r$ |

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

| Data source | w | d | X | | | Y | | | Z | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Design weights | Calibration weights | Calibration variables | | | Target variables | | | Response variables | | |
| Non–probability sample ($n$) | ? | ? | $x_{11}^n$ | $\cdots$ | $x_{1p}^n$ | $y_{11}^n$ | $\cdots$ | $y_{1p}^n$ | $z_{11}^n$ | $\cdots$ | $z_{1p}^n$ |
| | | | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| | | | $x_{n^n1}^n$ | $\cdots$ | $x_{n^np}^n$ | $y_{n^n1}^n$ | $\cdots$ | $y_{n^np}^n$ | $z_{n^n1}^n$ | $\cdots$ | $z_{n^np}^n$ |
| Calibration target data ($c$) | $w_1^c$ $\vdots$ $w_{n^c}^c$ | $d_1^c$ $\vdots$ $d_{n^c}^c$ | $x_{11}^c$ $\vdots$ $x_{n^c1}^c$ | $\cdots$ $\ddots$ $\cdots$ | $x_{1p}^c$ $\vdots$ $x_{n^cp}^c$ | ? | | | | | |
| Response–reference data ($r$) | $w_1^r$ $\vdots$ $w_{n^r}^r$ | $d_1^r$ $\vdots$ $d_{n^r}^r$ | | | | ? | | | $z_{11}^r$ $\vdots$ $z_{n^r1}^r$ | $\cdots$ $\ddots$ | $z_{1p}^r$ $\vdots$ $z_{n^rp}^r$ |

— Response model        — Calibration model        — Prediction model

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

| Data source | $w$ | $d$ | $X$ | | | $Y$ | | | $Z$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Design weights | Calibration weights | Calibration variables | | | Target variables | | | Response variables | | |
| Non–probability sample ($n$) | ? | ? | $x_{11}^n$ | $\cdots$ | $x_{1p}^n$ | $y_{11}^n$ | $\cdots$ | $y_{1p}^n$ | $z_{11}^n$ | $\cdots$ | $z_{1p}^n$ |
| | | | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| | | | $x_{n^n1}^n$ | $\cdots$ | $x_{n^np}^n$ | $y_{n^n1}^n$ | $\cdots$ | $y_{n^np}^n$ | $z_{n^n1}^n$ | $\cdots$ | $z_{n^np}^n$ |
| Calibration target data ($c$) | $w_1^c$ $\vdots$ $w_{n^c}^c$ | $d_1^c$ $\vdots$ $d_{n^r}^c$ | $x_{11}^c$ $\vdots$ $x_{n^c1}^c$ | $\cdots$ | $x_{1p}^c$ $\vdots$ $x_{n^cp}^c$ | | ? | | | | |
| Response–reference data ($r$) | $w_1^r$ $\vdots$ $w_{n^r}^r$ | $d_1^r$ $\vdots$ $d_{n^r}^r$ | | | | | ? | | $z_{11}^r$ $\vdots$ $z_{n^r1}^r$ | $\cdots$ $\ddots$ | $z_{1p}^r$ $\vdots$ $z_{n^rp}^r$ |

—— Response model ——————— Calibration model ——————— Prediction model

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

| Data source | w | d | X | | | Y | | | Z | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | Design weights | Calibration weights | Calibration variables | | | Target variables | | | Response variables | | |
| Non–probability sample ($n$) | ? | ? | $x_{11}^n$ $\cdots$ $x_{1p}^n$ $\vdots$ $\ddots$ $\vdots$ $x_{n^n1}^n$ $\cdots$ $x_{n^np}^n$ | | | $y_{11}^n$ $\cdots$ $y_{1p}^n$ $\vdots$ $\ddots$ $\vdots$ $y_{n^n1}^n$ $\cdots$ $y_{n^np}^n$ | | | $z_{11}^n$ $\cdots$ $z_{1p}^n$ $\vdots$ $\ddots$ $\vdots$ $z_{n^n1}^n$ $\cdots$ $z_{n^np}^n$ | | |
| Calibration target data ($c$) | $w_1^c$ $\vdots$ $w_{n^c}^c$ | $d_1^c$ $\vdots$ $d_{n^r}^c$ | $x_{11}^c$ $\cdots$ $x_{1p}^c$ $\vdots$ $\ddots$ $\vdots$ $x_{n^c1}^c$ $\cdots$ $x_{n^cp}^c$ | | | ? | | | | | |
| Response–reference data ($r$) | $w_1^r$ $\vdots$ $w_{n^r}^r$ | $d_1^r$ $\vdots$ $d_{n^r}^r$ | | | | ? | | | $z_{11}^r$ $\cdots$ $z_{1p}^r$ $\vdots$ $\ddots$ $\vdots$ $z_{n^r1}^r$ $\cdots$ $z_{n^rp}^r$ | | |

—— Response model          —— Calibration model          —— Prediction model

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

| Data source | Design weights $w$ | Calibration weights $d$ | Calibration variables $X$ | | | Target variables $Y$ | | | Response variables $Z$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Non–probability sample ($n$) | ? | ? | $x_{11}^n$ | $\cdots$ | $x_{1p}^n$ | $y_{11}^n$ | $\cdots$ | $y_{1p}^n$ | $z_{11}^n$ | $\cdots$ | $z_{1p}^n$ |
| | | | $\vdots$ | $\ddots$ | $\vdots$ | $\vdots$ | | $\vdots$ | $\vdots$ | $\ddots$ | $\vdots$ |
| | | | $x_{n^n1}^n$ | $\cdots$ | $x_{n^np}^n$ | $y_{n^n1}^n$ | $\cdots$ | $y_{n^np}^n$ | $z_{n^n1}^n$ | | $z_{n^np}^n$ |
| Calibration target data ($c$) | $w_1^c$ $\vdots$ $w_{n^c}^c$ | $d_1^c$ $\vdots$ $d_{n^r}^c$ | $x_{11}^c$ $\vdots$ $x_{n^c1}^c$ | $\cdots$ $\ddots$ $\cdots$ | $x_{1p}^c$ $\vdots$ $x_{n^cp}^c$ | ? | | | | | |
| Response–reference data ($r$) | $w_1^r$ $\vdots$ $w_{n^r}^r$ | $d_1^r$ $\vdots$ $d_{n^r}^r$ | | | | ? | | | $z_{11}^r$ $\vdots$ $z_{n^r1}^r$ | $\cdots$ $\ddots$ | $z_{1p}^r$ $\vdots$ $z_{n^rp}^r$ |

—— Response model        —— Calibration model        —— Prediction model

Statistical inference and quality
**Compensation methods**
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## Response model

▶ Modelling the response process:
  Approximation of unknown inclusion / participation probability

  $$\pi_i \;=\; P\left(i \in \mathcal{S}^n\right) \;=\; P\left(R_i = 1\right)$$

  where

  $$R_i \;=\; \begin{cases} 1, & i \in \mathcal{S}^n \\ 0, & i \in \mathcal{S}^r \end{cases}$$

▶ Set of observations $\mathcal{S}^r$ outside non–probability sample
  required to estimate model

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## Response model

▶ by logistic regression model on response–variables $\boldsymbol{Z}$ using model parameters $\boldsymbol{\omega}$:

$$\pi_i \approx \widehat{\pi}_i := \left(1 + \exp\left(-\boldsymbol{z}_{i.}^{\mathsf{T}} \boldsymbol{\omega}\right)\right)^{-1} \tag{1}$$

▶ Estimation via weighted binomial log-likelihood (pseudo-log-likelihood) of $\boldsymbol{R}$ given $\boldsymbol{\omega}$ and $\boldsymbol{Z}$

$$\begin{aligned}
\log\left(\mathcal{L}_{\mathrm{o}}\left(R|\boldsymbol{\omega}, \boldsymbol{Z}^n, \boldsymbol{Z}^r, \boldsymbol{w}^n, \boldsymbol{w}^r\right)\right) \\
= \sum_{i \in \mathcal{S}^n} w_i^n \cdot \log\left(\widehat{\pi}_i\right) + \sum_{i \in \mathcal{S}^r} w_i^r \cdot \log\left(1 - \widehat{\pi}_i\right)
\end{aligned}$$

▶ (Initial) weights $\boldsymbol{w}^n$ and $\boldsymbol{w}^r$ *might* be ones
(cf. Fuller, 2009; Rosenbaum und Rubin, 1983; Valliant und Dever, 2011)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Response model: Individual weighting

▶ Non–probability sampling treated as (single or additional) stage of random sampling with unknown probabilites

▶ Assumption: Participation is a random phenomenon

▶ $\widehat{\pi}_i$ is treated like inclusion probability in random sampling

$$\widetilde{w}_i := w_i^n \cdot \widehat{\pi}_i^{-1} \quad \text{for all} \quad i \tag{2}$$

(cf. e.g. Little, 1988; Valliant und Dever, 2011)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Response model: Grouped weighting

▶ Model (2): Possibly high variance of estimators

▶ Replacing $\widehat{\pi}_i$ by mean of similar observations

$$\widetilde{w}_i \;=\; w_i^n \cdot \left( \sum_{j \in g} 1 \right) \Big/ \left( \sum_{j \in g} \widehat{\pi}_j \right) \quad \text{for all} \quad i \in g$$

(3)

$g$: class including observation $i$

▶ Aim:

Reduce variability of weights

Less vulnerability to model misspecification

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## Response model: Grouped weighting: Post–stratification

▶ Weighted class proportions are calibrated to those of $\mathcal{S}^r$

▶ Replacing $\widehat{\pi}_i$ by post–stratification weight of propensity classes

$$\widetilde{w}_i \;=\; w_i^n \cdot \left( \sum_{j \in (g \cup \mathcal{S}^r)} w_j^n \right) \Big/ \left( \sum_{j \in (g \cup \mathcal{S}^n)} w_j^r \right) \quad \text{for all} \quad i \in g$$

(4)

$g$: class including observation $i$

(cf. e.g. Little, 1986; Rosenbaum und Rubin, 1983; Valliant und Dever, 2011)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibration model

▶ Find weights such that estimates meet known totals

$$\boldsymbol{\tau}\left(\boldsymbol{X}^{n}, \widetilde{\boldsymbol{w}} \circ \boldsymbol{d}\right) \;\stackrel{!}{=}\; \boldsymbol{\tau}\left(\boldsymbol{X}\right) \tag{5}$$

Different ways to achieve calibration constraints, e.g.

▶ Generalized regression estimator (GREG)
*(includes post–stratification)*

▶ Raking

(cf. Deville und Särndal, 1992; Deming und Stephan, 1940)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibration model: Targets' quality

**Known population totals as calibration targets**

▶ Exact compliance (often) reasonable

▶ Possibly high variance in weights / estimates if $\boldsymbol{X}$ includes many variables

**Estimated calibration targets**

▶ Commonly, high–quality random samples are used

▶ Subject to (survey–)errors as well

▶ Exact compliance less reasonable

▶ Inexact (relaxed) calibration is considered

(cf. Chang und Kott, 2008; Deville und Särndal, 1992; Deville et al., 1993)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## Calibration model: Relaxed constraints

▶ Exact compliance (equation (5)) is replaced by an adequate similarity:

$$\boldsymbol{\tau}\left(\boldsymbol{X}^n, \widetilde{\boldsymbol{w}}\right) \ \overset{!}{=} \ \boldsymbol{\tau}\left(\boldsymbol{X}^c, \boldsymbol{w}^c\right) \circ \boldsymbol{\epsilon} \qquad (6)$$

▶ $\boldsymbol{\epsilon}$ is a multiplicative error vector, determining the relation

$$\epsilon_k = \frac{\boldsymbol{\tau}\left(\boldsymbol{x}^n_{\cdot k}, \widetilde{\boldsymbol{w}}\right)}{\boldsymbol{\tau}\left(\boldsymbol{x}^c_{\cdot k}, \boldsymbol{w}^c\right)}$$

of the $k$–th total:

$\epsilon_k \ < \ 1$:   below target

$\epsilon_k \ = \ 1$:   exact compliance

$\epsilon_k \ > \ 1$:   above target

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibration model: Relaxed constraints

Guggemos und Tillé (2010): 'penalized Calibration' resembles GREG:

$$\underset{\boldsymbol{\omega},\boldsymbol{\epsilon}}{\operatorname{argmin}} \left( \sum_{j=1}^{n^n} w_j^n \cdot \frac{(1-d_j)^2}{2} + \sum_{k=1}^{p} \mathsf{v}_k \cdot \frac{(1-\epsilon_k)^2}{2} \right)$$

$$\text{s. t.} \quad \boldsymbol{\tau}\left(\boldsymbol{X}^n, \widetilde{\boldsymbol{w}}\right) \; \overset{!}{=} \; \boldsymbol{\tau}\left(\boldsymbol{X}^c, \boldsymbol{w}^c\right) \circ \boldsymbol{\epsilon} \tag{7}$$

$$\mathsf{L}_{\epsilon_k} \; \leq \; \epsilon_k \; \leq \; \mathsf{U}_{\epsilon_k} \quad \text{for all} \quad k = 1, \ldots, p$$

▶ Each $\epsilon_k$ is either fixed to 1 or unconstrained

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibration model: Gelman bounds

▶ Limiting range / variation of weights: '*Gelman-bounds*':

$$\frac{\text{Max}\,(\widetilde{\boldsymbol{w}})}{\text{Min}\,(\widetilde{\boldsymbol{w}})} \tag{8}$$

▶ Additional *boundary constraints* are introduced:

$$\text{L}_{\boldsymbol{d}} \;\leq\; d_j \;\leq\; \text{U}_{\boldsymbol{d}} \qquad \text{for all} \quad j \;=\; 1, \ldots, n^n \tag{9}$$

▶ $\text{L}_{\boldsymbol{d}}$ and $\text{U}_{\boldsymbol{d}}$ are global lower and upper bounds for $\boldsymbol{d}$

(cf. Gelman, 2007; Meng et al., 2009; Münnich et al., 2012a)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibration model: Extensions

**Münnich et al. (2012c):**

▶ Arbitrary Box–constraints

$L_{\widetilde{w}_j}$ and $U_{\widetilde{w}_j}$ for weights / Gelman factor and

$L_{\epsilon_k}$ and $U_{\epsilon_k}$ for totals

**Münnich et al. (2012b):**

▶ Computational effective implementation

▶ Using duality–approach to

▶ Determine weights by Lagrange multipliers, thus reducing dimensions to number of constraints

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibration model: Functional form approach

▶ Different distance functions often lead to very similar results

▶ Advantage of distance functions is questioned

▶ *Functional form approach* / *instrument vector approach*

▶ Correction depending on *instrument variables* $\boldsymbol{Z}$ and parameters $\boldsymbol{\omega}$:

$$\boldsymbol{d} \ := \ \mathbf{d}_{\mathrm{o}}\left(\boldsymbol{\omega}, \boldsymbol{Z}^{n}\right) \tag{10}$$

▶ $\boldsymbol{Z}$ does not (necessarily) coincide with $\boldsymbol{X}$

(cf. Estevao und Särndal, 2000, 2006; Folsom und Singh, 2000)

Statistical inference and quality
**Compensation methods**
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibrated response model

▶ *Generalized raking model*

$$\mathbf{d}_r\left(\boldsymbol{\omega}, \boldsymbol{Z}^n\right) = \exp\left(\boldsymbol{Z}^n\boldsymbol{\omega}\right) \tag{11}$$

▶ *Logit model* (inverse propensity score, cf. equation (1))

$$\mathbf{d}_l\left(\boldsymbol{\omega}, \boldsymbol{Z}^n\right) = 1 + \exp\left(-\boldsymbol{Z}^n\boldsymbol{\omega}\right) \tag{12}$$

*Calibrated propensity weights*, but possibly differing parameter estimation:

▶ If $\dim\left(\boldsymbol{X}^n\right) = \dim\left(\boldsymbol{Z}^n\right)$: $\boldsymbol{\omega}$ determined from constraints alone
▶ Otherwise, distance functions still needed

(cf. Folsom und Singh, 2000; Kott, 2003, 2006)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Calibrated response model: Distance functions & relaxed constraints

▶ Chang und Kott (2008) propose minimizing the distance to calibration targets for raking– or logit–model (11) and (12)

$$\operatorname*{argmin}_{\boldsymbol{\omega}} \left( \sum_{k=1}^{p} \mathsf{v}_k \cdot \frac{(1-\epsilon_k)^2}{2} \right) \tag{13}$$

$$\text{s.t.} \quad \boldsymbol{\tau}\left(\boldsymbol{X}^n, \widetilde{\boldsymbol{w}}\right) \stackrel{!}{=} \boldsymbol{\tau}\left(\boldsymbol{X}^c, \boldsymbol{w}^c\right) \circ \boldsymbol{\epsilon}$$

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## Prediction model: Linear Regression

▶ Regression equation

$$\widehat{y}_{il}^n = \beta_0 + \sum_{j=1}^{p} x_{ij}^n \cdot \beta_j \qquad . \qquad (14)$$

▶ $\boldsymbol{x}_{\cdot 0}^n$: intercept column,

▶ Parameters determined by least squares

$$\boldsymbol{\beta} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left( (\boldsymbol{E}^n)^{\mathsf{T}} \operatorname{diag}(\boldsymbol{w}^n) \, \boldsymbol{E}^n \right)$$

$$= \left( (\boldsymbol{X}^n)^{\mathsf{T}} \operatorname{diag}(\boldsymbol{w}^n) \, \boldsymbol{X}^n \right)^{-1} (\boldsymbol{X}^n)^{\mathsf{T}} \operatorname{diag}(\boldsymbol{w}^n) \, \boldsymbol{Y}^n \qquad (15)$$

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## Prediction model: Support vector machine

▶ Regression equation

$$\widehat{y}_{il}^n = \beta_0 + \sum_{j=1}^p x_{ij}^n \cdot \beta_j \qquad (16)$$

▶ Parameters determined by

$$\boldsymbol{\beta} = \operatorname*{argmin}_{\boldsymbol{\beta}} \left( \frac{1}{2} \cdot \sum_{j=1}^p \beta_j^2 + C \cdot \sum_{i=1}^{n^n} \xi_i + C \cdot \sum_{i=1}^{n^n} \xi_i \right) \qquad (17)$$

$$\text{s.t.} \qquad \xi_i, \; \xi_i \; \geq \; 0$$

$$|y_{il}^n - \widehat{y}_{il}^n| \; \leq \; e + \xi_i$$

▶ Slack-variables $\xi_i$, $\xi_i^*$ for violation of the

▶ *Maximum* distance $e$ of points to the regression line

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Further predictive methods

From the methodological point of view, any prediction model can be applied. However, they are relying on strong assumptions, and likely additional weighting, benchmarking or alignment methods might be considered.

A very detailed reading is:

Simon Lenau (2023): Statistical and Machine Learning Methods for Handling Selectivity in Non-Probability Samples. PhD dissertation, Trier University.

Further readings: InGRID deliverable at
https://www.inclusivegrowth.eu/project-output
and the references therein.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Graphical representation of missingness patterns

Two data sources. '✓' and '?' indicate observed and missing data

|  |  | X | Y |
|---|---|---|---|
| Probability Sample | 1 | ✓ | ? |
|  | ⋮ | ⋮ | ⋮ |
|  | n | ✓ | ? |
| Big Data Sample | 1 | ✓ | ✓ |
|  | ⋮ | ⋮ | ⋮ |
|  | $N_B$ | ✓ | ✓ |

See Yang and Kim (2018)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Combining information from individual records

In some cases, data records for individuals can be merged from different sources. There are different methods applicable:

▶ Record linkage (possibly legal constraints)

▶ Statistical matching

▶ Multiple imputation

▶ Mass imputation

Again, model assumptions have to be met (e.g. CIA), some of which can be hardly verified.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Comparison of income class distributions (Wage indicator vs. Microcensus)

Statistical inference and quality
Compensation methods
**Web surveys and big data**
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Absolute difference of income class frequencies (Wage indicator vs. Microcensus)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Relative difference of income class frequencies (Wage indicator vs. Microcensus)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Simulation setup

▶ Population: Subset of the $\mathrm{AMELIA}$ universe with $N = 20\,000$

▶ $R = 1\,000$ Poisson–samples of size $n = 1\,000$ (5%)

  Participation probabilities $\pi$

▶ Fixed correlations with variables

  gender

  isced

  bas

  age

▶ Varying correlation with income variable inc

▶ $\pi$ assumed to be unknown

▶ Estimation with presented methods

  (cf. Burgard et al., 2017)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Estimated mean income

$\mu$ (inc)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Estimated correlation between income and household–income

$\rho\,(inc, hhinc)$



Calibration targets: Population

Calibration & prediction model variables: gender isced bas age

Response model variables: gender isced bas age hhinc
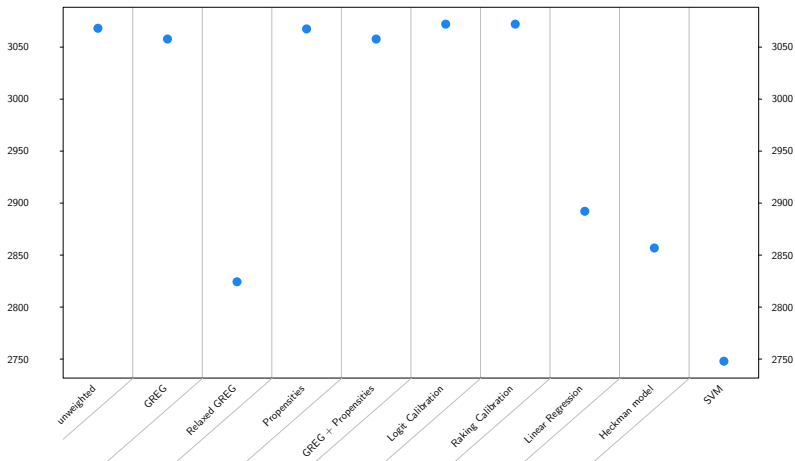
Statistical inference and quality
Compensation methods
**Web surveys and big data**
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Estimates from the WageIndicator Survey:



Mean income

Statistical inference and quality
Compensation methods
**Web surveys and big data**
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Estimates from the WageIndicator Survey:

Mean household income

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Two cases of use of Big Data

▶ Twitter sentiment analysis
▶ Quality measure of Satellite images

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Twitter Sentiment Analysis

▶ Sentiment analysis is a natural language processing (NLP) technique used to determine whether data is positive or negative.

▶ Tweets collected using #JoeBiden, #DonaldTrump, #Biden, #Trump with the Twitter API

▶ 38,432,811 tweets were collected, employing streaming Tweepy API across the United States between 28 September 2020, and 20 November 2020.

See Chaudhry et al. (2021)

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Figure: Twitter dataset

| | created_at | tweet_id | tweet | likes | retweet_count | source | user_id | user_name |
|---|---|---|---|---|---|---|---|---|
| 0 | 2020-10-15 00:00:01 | 1.316529e+18 | #Elecciones2020 \| En #Florida: #JoeBiden dice ... | 0.0 | 0.0 | TweetDeck | 3.606665e+08 | El Sol Latino News |
| 1 | 2020-10-15 00:00:18 | 1.316529e+18 | #HunterBiden #HunterBidenEmails #JoeBiden #Joe... | 0.0 | 0.0 | Twitter for iPad | 8.099044e+08 | Cheri A. us |
| 2 | 2020-10-15 00:00:20 | 1.316529e+18 | @IslandGirlPRV @BradBeauregardJ @MeidasTouch T... | 0.0 | 0.0 | Twitter Web App | 3.494182e+09 | Flag Waver |
| 3 | 2020-10-15 00:00:21 | 1.316529e+18 | @chrislongview Watching and setting dvr. Let's... | 0.0 | 0.0 | Twitter for iPhone | 8.242596e+17 | Michelle Ferg |
| 4 | 2020-10-15 00:00:22 | 1.316529e+18 | #censorship #HunterBiden #Biden #BidenEmails #... | 1.0 | 0.0 | Twitter Web App | 1.032807e+18 | the Gold State |

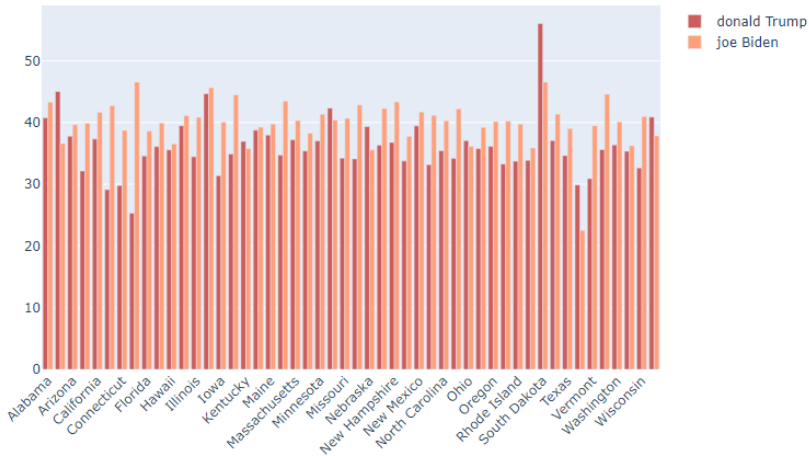Statistical inference and quality
Compensation methods
**Web surveys and big data**
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Positivity proportion per state

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Figure: Sentiment analysis and real poll results for Joe Biden

| | Maine | California | New York | Arkansas | Idaho |
|---|---|---|---|---|---|
| Sentiments | 64.3% | 64.2% | 61% | 62.1% | 60.8% |
| Margin of Victory | 9% | 30% | 23.2% | $-27.6\%$ | $-30.7\%$ |

▶ Results in Maine, California, and New York aligned with the Twitter sentiment analysis

▶ Arkansas and Idaho although had positive Twitter sentiment, however, had opposing responses in the real poll.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Using satellite images

▶ Satellite-based recordings can have global coverage
▶ They are considered objective measurements
▶ Satellites take multiple recordings
▶ Some data are available for free
▶ Can act as a projection space to relate multiple measurements together with environmental information.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Improving Forest Inventory

Julian Wagner and Ralf Münnich et al. (2017) used satellite data to improve the Rhineland-Palatinate forest inventory.

▶ The standard German forest inventory (GNFI) selects some forest areas and collects high-quality information on about 150 variables.

▶ For local, small regions the sample size is insufficient.

▶ Airborn Laser scanning (ALS) data from 2002-2013 were available and combined with topographic maps.

▶ The forest canopy height minus the ground elevation results in normalized surface models.

▶ This measure of vegetation height is used as a proxy variable to improve the GNFI survey data in a small area estimation

▶ Problem: The canopy height is related non-linear to biomass.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

Given:

▶ $U$ a population of size $N$ positions in $D$ areas with populations sizes $N_d$

▶ A sample of $n$ partitioned into subsamples $S_d$ of size $n_d$.

▶ Estimating the small area mean $\theta$ for each area with:

$$\theta_d = \mu_{y,d} := \frac{1}{N_d} \sum_{i \in U_d} y_i$$

   $n_d$ might be to small for reliable estimates of $\mu_{y,d}$ for each area.

▶ Canopy heights might be a good proxy to improve on the local estimates $\mu_{y,d}$

▶ Standard small area approaches are parametric model-based, implying a linear relationship of the dependent variable and predictors in the form of a random effect (see Battese et al. 1998).

$$y_{i_d} = \tilde{x}_{i,d}^T + u_d + e_{i,d}$$

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

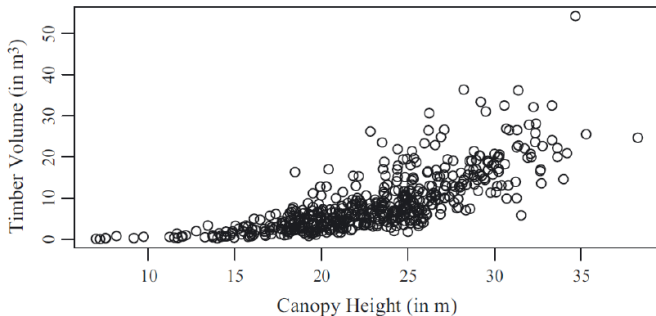Tree height and timber volume are not linear:



Figure 1 from: Wagner et al. (2017): Non-parametric small area models using shape-constrained penalized B-splines. in Royal Statistical Society, p.1089–1109

$\rightarrow$ non-parametric, B-spline based estimation.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Data quality in NPS and big data

▶ Classical concepts of data quality or the total survey error do not really meet the needs for big data

▶ Often non-response measures are used (e.g. R indicators, cf. Shlomo, Schouten, and others)

▶ Fitness for use (user oriented, cf. Wang and Strong, 1996)

▶ Approaches are in development (e.g. Biemer, 2017)

▶ But inference (NPS and Big Data) is still an issue, and for big data additionally the complexity

Discussion and the references:

Münnich/Articus (2022): Big Data und Qualität - ist viel gleich gut? 85-99. In: B. Wawrzyniak/M. Herter (Eds.): Neue Dimensionen in Data Science, Berlin.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Summary

All approaches require auxiliary data:

### Response models:

▶ Reference data with good response–predictors

### Calibration models:

▶ Calibration targets, correlating with variables of interest

### Prediction models:

▶ Totals / means (*linear* models), microdata (*non–linear* models)

▶ Good predictors for **every variable of interest**

**Editing / Data cleaning of utmost importance**

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Issues with NPS and Big Data

▶ There is no general strategy but case-specific
▶ Weighting helps but in special cases, model-based methods might be better – but how to know?
▶ NPS have big advantages, when time matters, e.g. changes right before elections or special events (influential information)
▶ Big data can help a lot
    ▶ mobile data: traffic control, or road use
    ▶ satellite data: forest inventory, urban patterns, (independant) SDG indicators
▶ BUT: whenever we aim getting important information (e.g. for budget transfers), use proper high-quality survey data!

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Issues with NPS and Big Data

▶ There is no general strategy but case-specific

▶ Weighting helps but in special cases, model-based methods might be better – but how to know?

▶ NPS have big advantages, when time matters, e.g. changes right before elections or special events (influential information)

▶ Big data can help a lot
  ▶ mobile data: traffic control, or road use
  ▶ satellite data: forest inventory, urban patterns, (independant) SDG indicators

▶ BUT: whenever we aim getting important information (e.g. for budget transfers), use proper high-quality survey data!

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# Thanks for your attention!

And thanks to the InGRID Research Infrastructure who

financially supported the research on the wage indicator survey

(https://www.inclusivegrowth.eu) as well as Simon

Lenau and Abrar Ahmed for their support.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# References I

Burgard, J. P., Kolb, J.-P., Merkle, H. und Münnich, R. (2017):
   *Synthetic data for open and reproducible methodological
   research in social sciences and official statistics*. In: AStA
   Wirtschafts- und Sozialstatistisches Archiv, 11 (3), S. 233–244.

Chang, T. und Kott, P. S. (2008): *Using calibration weighting to
   adjust for nonresponse under a plausible model*. In: Biometrika,
   95 (3), S. 555–571.

Chaudhry, H. N., Javed, Y., Kulsoom, F., Mehmood, Z., Khan,
   Z. I., Shoaib, U. und Janjua, S. H. (2021): *Sentiment analysis of
   before and after elections: Twitter data of US election 2020*. In:
   Electronics, 10 (17), S. 2082.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# References II

Deming, W. E. und Stephan, F. F. (1940): *On a Least Squares Adjustment of a Sampled Frequency Table When the Expected Marginal Totals are Known*. In: The Annals of Mathematical Statistics, 11 (4), S. 427–444.

Deville, J.-C. und Särndal, C.-E. (1992): *Calibration estimators in survey sampling*. In: Journal of the American statistical Association, 87 (418), S. 376–382.

Deville, J.-C., Särndal, C.-E. und Sautory, O. (1993): *Generalized raking procedures in survey sampling*. In: Journal of the American statistical Association, 88 (423), S. 1013–1020.

Estevao, V. M. und Särndal, C.-E. (2000): *A functional form approach to calibration*. In: Journal of Official Statistics, 16 (4), S. 379–399.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# References III

Estevao, V. M. und Särndal, C.-E. (2006): *Survey estimates by calibration on complex auxiliary information*. In: International Statistical Review, 74 (2), S. 127–147.

Folsom, R. E. und Singh, A. C. (2000): *The generalized exponential model for sampling weight calibration for extreme values, nonresponse, and poststratification*. In: Proceedings of the Survey Research Methods Section, American Statistical Association (2000), S. 598–603.

Fuller, W. A. (2009): Sampling Statistics. Wiley series in survey methodology.

Gelman, A. (2007): *Struggles with survey weighting and regression modeling*. In: Statistical Science, 22 (2), S. 153–164.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# References IV

Guggemos, F. und Tillé, Y. (2010): *Penalized calibration in survey sampling: Design-based estimation assisted by mixed models*. In: Journal of statistical planning and inference, 140 (11), S. 3199–3212.

Kott, P. S. (2003): *A practical use for instrumental-variable calibration*. In: Journal of Official Statistics, 19 (3), S. 265.

Kott, P. S. (2006): *Using calibration weighting to adjust for nonresponse and coverage errors*. In: Survey Methodology, 32 (2), S. 133–142.

Little, R. J. (1986): *Survey nonresponse adjustments for estimates of means*. In: International Statistical Review, 54 (2), S. 139–157.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

## References V

Little, R. J. (1988): *Missing-data adjustments in large surveys*. In: Journal of Business & Economic Statistics, 6 (3), S. 287–296.

Meng, X., Duan, N., Chen, C. und Alegria, M. (2009): *Power-shrinkage: An alternative method for dealing with excessive weights*. In: Proceedings of the joint statistical meetings 2009.

Münnich, R., Gabler, S., Ganninger, M., Burgard, J. P. und Kolb, J.-P. (2012a): *Stichprobenoptimierung und Schätzung im Zensus 2011*. Statistik und Wissenschaft Bd. 21.

Münnich, R., Sachs, E. W. und Wagner, M. (2012b): *Calibration of estimator-weights via semismooth Newton method*. In: Journal of Global Optimization, 52 (3), S. 471–485.

Statistical inference and quality
Compensation methods
Web surveys and big data
Conclusion and outlook

Lehrstuhl für Wirtschafts- und Sozialstatistik

# References VI

Münnich, R. T., Sachs, E. W. und Wagner, M. (2012c): *Calibration of estimator-weights via semismooth Newton method*. In: Journal of Global Optimization, 52 (3), S. 471–485.

Rosenbaum, P. R. und Rubin, D. B. (1983): *The central role of the propensity score in observational studies for causal effects*. In: Biometrika, 70 (1), S. 41–55.

Valliant, R. und Dever, J. A. (2011): *Estimating propensity adjustments for volunteer web surveys*. In: Sociological Methods & Research, 40 (1), S. 105–137.