

# Data Quality: Total Survey Error (TSE)

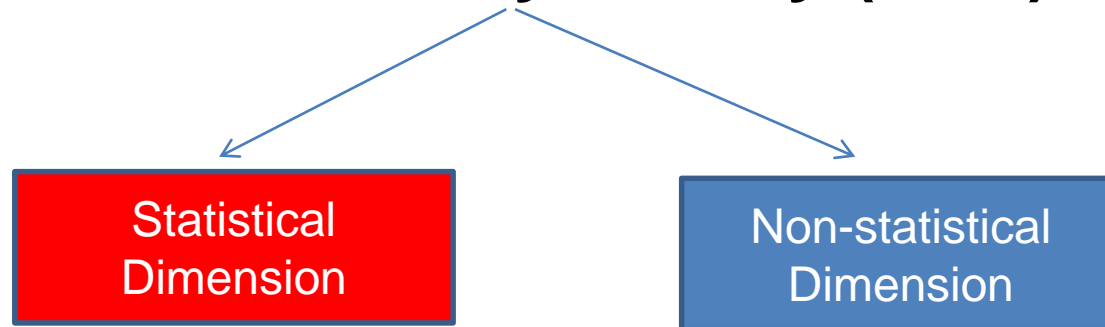
Olga Maslovskaya  
University of Southampton

## Survey Data

- Vast amounts of survey data are collected for many purposes, including governmental information, public opinion and election surveys, advertising and market research as well as scientific research
- Survey data underlie many public policy and business decisions
- Good quality data reduces the risk of poor policies and decisions and is of crucial importance

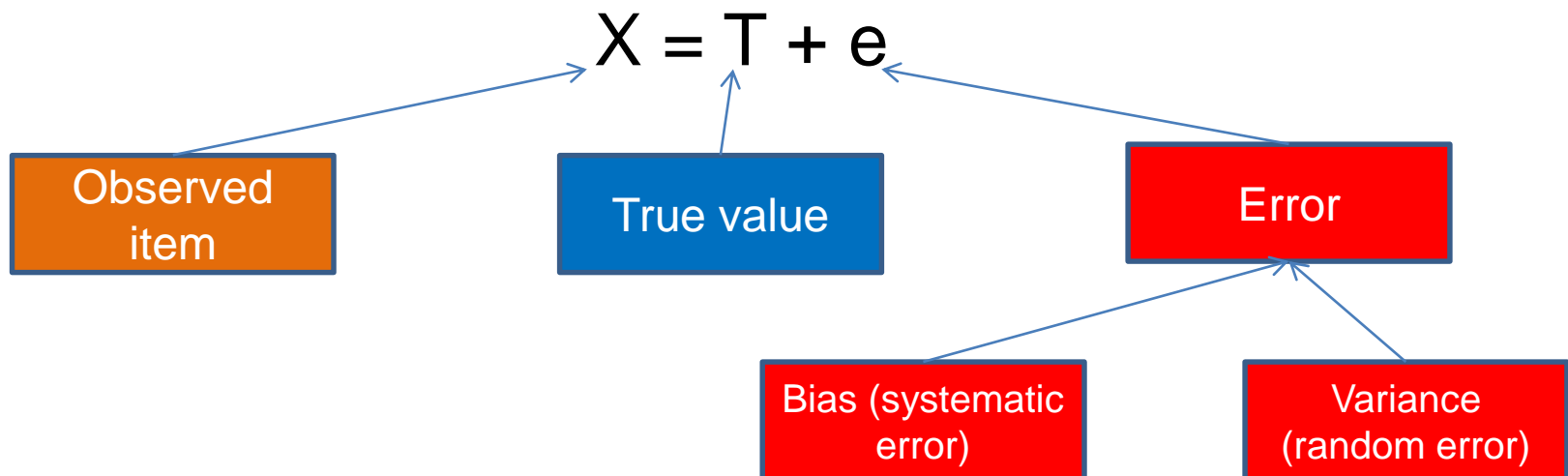
# Total Survey Quality (TSQ)

## Total Survey Quality (TSQ)



## TSQ: Quality Dimensions –Statistical

- **Accuracy** of estimates is the difference between the estimate and the true parameter value
- **Accuracy** is the larger concept of TSQ



## Total Survey Error (TSE) (I)

- TSE concept was developed by Robert Groves (1989) in book on Survey Errors and Survey Costs
- Survey estimates are derived from complex survey data, published estimates may differ from their true parameter values due to survey errors
- Total Survey Error is the difference between a population mean, total, or other population parameter and the estimate of the parameter based on the sample survey (or census) (Biemer and Lyberg, 2003)

## Total Survey Error (TSE) (2)

- Survey error is any error arising from the survey process that contributed to the deviation of an estimate from its true parameter value (Biemer, 2016)
- Survey error diminishes the accuracy of inferences derived from the survey
- TSE is the accumulation of **all** errors that may arise in the design, collection, processing, and analysis of survey data (Biemer, 2016)

## TSE framework (I)

- Set of principles, methods and processes that minimise TSE within the budget allocated for accuracy, timing and other constraints
- Non-statistical dimensions of TSQ can be viewed as constraints – timeliness and comparability constrain the design; accessibility, relevance and completeness constrain the budget (Biemer 2017)

## TSE framework (2)

TSE paradigm provides principles that guide stages of survey process:

- Survey design
- Implementation
  - Data collection
  - Data processing
  - Estimation
- Quality evaluation
- Data analysis

Each stage of survey process provides opportunities for errors which add up to TSE



## TSE

**TSE= sampling errors + non-sampling errors**

### Survey errors:

- **Sampling errors** – can be computed for probability samples and are due to selecting a sample instead of the entire population
- **Non-sampling errors** (including measurement error – cannot be formally estimated but can be improved by interviewing procedures and question wordings etc.) - are errors due to mistakes or system deficiencies, also from incomplete responses to the survey or its questions, etc.
- In many cases non-sampling error can be much more damaging than sampling error to estimates from surveys

# Sources of Sampling Error

- **Sampling scheme**
  - Stratification
  - Clustering
  - Selection probabilities
  - Sampling phases
- **Sample size**
  - Overall sample size
  - Effective sample size
  - Sample size allocation
- **Estimator choice**
  - Simple
  - Use of auxiliary information
  - Model-based
  - Model-assisted

## Components of Non-sampling Error

1. Specification error
2. Frame error
3. Nonresponse error
4. Measurement error
5. Processing error
6. Modelling/Estimation error

Biemer (2017)

## Specification Error

- Refers to a question on the questionnaire
- Occurs when the concept implied by the survey question and the concept that should be measured in the survey differ (Biemer and Lyberg, 2003)

## Frame Error

- Arises from construction of the sampling frame for the survey
- The sampling frame might have erroneous omissions, duplicates or erroneous inclusions

## Nonresponse Error

- **Unit nonresponse** occurs when a sample unit (individual, household or organisation) does not respond to any part of the questionnaire,
- **Item nonresponse** occurs when the questionnaire is only partially completed and some items are not answered
- **Incomplete response** occurs when the response to open-ended question is incomplete or very short and inadequate
- **Panel attrition** occurs when a sample unit is lost over the period of a longitudinal study

## Measurement error

- Measurement errors pose a serious limitation to the validity and usefulness of the data collected
- Most damaging source of error
- Having excellent samples representative of the target population, high response rates, complete data, etc. does us little good if our measurement instruments evoke responses that are fraught with error
- Without reliable measurements, analysis of data hardly make any sense

# Key components of measurement error

- Respondents
  - May deliberately or unintentionally provide incorrect information
    - Response style behaviours
    - Satisficing (less efforts to provide optimal responses)
- Interviewers - enumerators
  - May falsify data
  - May inappropriately influence responses
  - May have negative impact on responses to sensitive questions
  - May record responses incorrectly
  - May fail to comply with the survey protocol
- Questionnaire - design
  - Bad design
    - Ambiguous questions
    - Confusing instructions
    - Unclear terms
- Mode of administration
  - Online mode
    - Non-optimised questionnaire for smartphones



## Processing Error

Contributes to measurement error

- Occurs during data processing stage
  - Errors in data editing
  - Errors in data entry
  - Errors in coding
  - Errors in outlier editing
  - Errors in assignment of survey weights
  - Errors in non-response imputing

## Modelling and Estimation Error

Occurs during data analysis stage (modelling)

- Errors in weight adjustments,
- Errors in imputation,
- Errors in modelling process and in models

## Types of Errors

- **Systematic Error – *bias*** - errors that tend to agree – results in biased estimates (strengthen the relations between variables, leading to false conclusions) – e.g. response styles or other stable behaviours - bias the results, distorting the mean value on variables – does not cancel out
- **Random Error – *variance*** - errors that tend to disagree (unintended mistakes made by respondents) – affects the variance of estimates (may weaken the relations between variables), vary from case to case but are expected to cancel out

## Mean Squared Error (MSE)

- *Total survey error* (TSE) is a term that is used to refer to all sources of bias (systematic error) and variance (random error) that may affect accuracy of survey data.
- MSE is the sum of the total bias squared plus the variance components for all the various sources of error in the survey design.
- MSE – metric for measuring TSE
- MSE cannot be calculated directly but useful conceptually to consider how large the different components can be and how much they add to the total survey error
- Hypothetical but great guide for optimal survey designs

## MSE

- Survey design goal is to **minimise** the “mean squared error” (MSE)
- When other designs are similar on other quality dimensions, the optimal design is the one achieving the smallest mean squared error
- Working to reduce the measurement error on one set of questions could increase the error for a different set of questions in the same survey
- Also, reducing one error could increase another error in the survey

## Survey designers face the following questions:

- Where should additional resources be directed to generate the greatest improvement to data quality: extensive interviewer training for nonresponse reduction, greater nonresponse follow up intensity, or by offering larger incentives to sample members to encourage participation?
- Should a more expensive data collection mode be used, even if the sample size must be reduced significantly to stay within budget?

## TSE in Practice (I)

- Realistic scenario is to work on continuous improvement of various survey processes so that biases and unwanted variations are gradually reduced
  - Redesign of surveys if needed
  - Non-response bias reduction through responsive and adaptive survey designs
  - Data quality indicators application in data analysis
- Idea is to minimize all these error sources
- Minimizing all of these errors would require an unlimited budget (impossible)
- Cost-benefit trade-offs are needed to decide which errors to minimize

## TSE in Practice (2)

Decisions are needed:

- To ignore some errors
- To measure and to control/adjust for some (data analysis stage: complex designs, measurement errors, missing data, sampling errors)



## Conclusions

- Data accuracy is of crucial importance
- Single score or measure of data quality is not available
- Cost-benefit trade-offs to minimise different errors depending on survey aims
- TSE framework was developed and adopted
- TSE helps keeping data quality standards high and in line with survey aims under financial constraints

## References

- **Biemer** (2010) Total survey error: Design, implementation, and evaluation. *Public Opinion Quarterly*, 74(5): 817-848.
- **Biemer** (2016) Total Survey Error Paradigm: Theory and Practice. In *The Sage handbook of survey methodology* by Wolf, Joye, Smith and Fu. London: SAGE publications.
- **Biemer** (2017) Total survey error: A Framework for censuses and surveys. *Presentation* at the University of Southampton.
- **Biemer and Lyberg** (2003) *Introduction to survey quality*. New York: John Wiley & Sons.
- **Groves and Heeringa** (2006) Responsive design for household surveys: Tools for actively controlling survey errors and costs. *Journal of the Royal Statistical Society Series A*, 169 (3): 439-457.
- **Lynn** (2004) Editorial: Measuring and communicating survey quality. *Journal of the Royal Statistical Society Series A*, 167 (4): 575-578.
- **Lyberg and Weisberg** (2016) *The SAGE handbook of survey methodology*. London: SAGE publications.
- **Schouten et al.** (2013) Optimizing quality of response through adaptive survey designs. *Survey Methodology*, 39 (1): 29-39.
- **Weisberg** (2005) *The total survey error approach*. Chicago: University of Chicago Press.