Dealing with non-response: Introduction to imputation

Aksel Thomsen & Peter Tibert Stoltze Statistics Denmark

Dealing with non-response: Overview

- Defining non-response
- Assessment of non-response
- The problem with non-response
- Imputation as a way to deal with non-response
- Types of imputation
- Added uncertainty by applying imputation







- Assume a data file containing the results of a questionnaire with *p* questions (partially) answered by *n* subjects
- The data are organized in a $n \times p$ data frame*

$$Y = \begin{bmatrix} Y_{11} & \cdots & Y_{1p} \\ \vdots & \ddots & \vdots \\ Y_{n1} & \cdots & Y_{np} \end{bmatrix}$$

- This means that the answer from subject *i* to question *j* is Y_{ij}
- We call this a *data frame* because the columns need not be of the same type if all columns are numeric we have a matrix

• We now define a matrix *R* of same dimensions as *Y* to indicate non-response:

$$R_{ij} = \begin{cases} 1 & if \ Y_{ij} \in V_j \\ 0 & else \end{cases}$$

- Here we define the case of response to mean that a specific answer X_{ij} belongs to the V_j (the set of admissible values for question j)
- Hence, the matrix *R* becomes filled with 0's and 1's allowing us to start assessing the amount of non-response
- Note: This seems like a simple task, but it is not



- With non-response we move from *n* sampled units to *m* responding units
- If the mechanism from sample to response is truly random, then we can estimate the response propensity as

$$\theta = \frac{m}{n}$$

• The combined probability of both being sampled and responding (assuming these are actually independent) is simply

$$p = \pi\theta = \frac{n}{N}\frac{m}{n} = \frac{m}{N}$$



- However, θ is almost never a constant value, since response propensity almost always varies with our variables of interest
- In the 1970's Rubin suggested a more formal approach to this distinguishing between
 - Missing complete at random (MCAR) where in fact θ is the same for all cases
 - Missing (conditionally) at random (MAR) where θ can be constant within groups of observed data
 - Missing not at random (MNAR) where θ is individual
- Simple methods like Mean Imputation only works for MCAR, and they will give biased results for MAR and MNAR

Dealing with partial non-response

- Simplest solution: Dismiss all rows with partial non-response, i.e. reduce data set to *complete cases*
 - Simple but can potentially mean that too much information is lost
- A better solution is often to apply *imputation* to suitable variables
- Probably not all variables should be subject to imputation

Dealing with item non-response

- If a subject has not answered at all, we will only in rare circumstances do imputation
- Weighting is normally the best solution when no individual response is particularly important
- In the case of business statistics and large enterprises, it can make sense to do imputation
 - In this case we will resort to "expert imputation" which is a very manually driven process based on a combination of whatever knowledge is available – this process should only be applied to very few cases!

Types of imputation

- Donor based or model based
 - Duplicate values from a neighbor or construct from statistical/ML model
- Stochastic or deterministic
 - Get exact same result each time or introduce a bit of variation
- Hot-deck or cold-deck
 - Using information from current round of survey or from previous ones

Examples of imputation methods

- Mean imputation
- Last Value Carried Forward
- Multiple Linear regression
- Multinomial logistic regression
- Random forest





