

# Clinical prediction models

Session 11  
Dealing with missing data

Pedro E A A do Brasil  
pedro.brasil@ini.fiocruz.br  
2018









# Rationale

**Table 7.1** Hypothetical missing data pattern: 250 subjects have partially complete data (missing data indicated with . ), and 250 have fully complete data (indicated with X)

ID	X1	X2	X3	X4	X5	Y
1	.	X	X	X	X	X
...	.	X	X	X	X	X
50	.	X	X	X	X	X
51	X	.	X	X	X	X
...	X	.	X	X	X	X
100	X	.	X	X	X	X
101	X	X	.	X	X	X
...	X	X	.	X	X	X
150	X	X	.	X	X	X
151	X	X	X	.	X	X
...	X	X	X	.	X	X
200	X	X	X	.	X	X
201	X	X	X	X	.	X
...	X	X	X	X	.	X
250	X	X	X	X	.	X
251	X	X	X	X	X	X
...	X	X	X	X	X	X
...	X	X	X	X	X	X
500	X	X	X	X	X	X
Total	450	450	450	450	450	500

Steyerberg. Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating. Springer in 2009.



# Rationale

- Basics
  - For example, one may wish to compare nested models, or adjust analysis and have an idea of the adjusted effect from univariable to multivariable
  - In two models conducted with missing data, it is then impossible to infer whether differences in odds ratios, p values or R<sup>2</sup> arose because of true differences, because of correlation between the predictors or because of a selection of subjects due to missing values

# Missing mechanisms

- Depending of the imputation strategy, the mechanism is not that relevant.
- In health data the mechanism is usually not at random.

**Table 7.2** Three types of missing data mechanisms

Label	Missing mechanism	Description
MCAR	Missing completely at random	Administrative errors, accidents
MAR	Missing at random	Missingness related to known patient characteristics, time or place (“MAR on $x$ ”), or to the outcome (“MAR on $y$ ”)
MNAR	Missing not at random	Missingness related to the value of the predictor, or to characteristics not available in the analysis

Steyerberg. Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating. Springer in 2009.

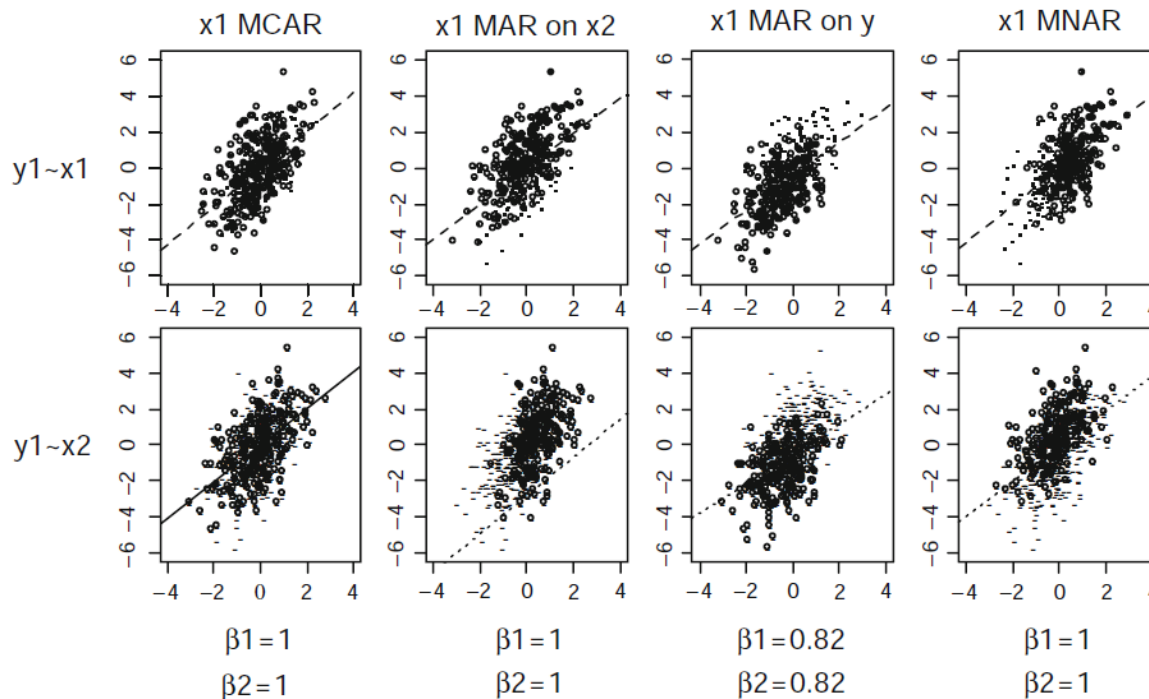


# Examples of bias

- A correlation between missingness of a predictor and the outcome poses a serious problem in predictive modelling.
- If an association between missingness of predictors  $X$  and outcome  $Y$  is noted in a prospective study, the explanation must be through other predictors.
- MAR on  $y$  for only one predictor is sufficient to bias coefficients of all predictors.

Steyerberg. Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating. Springer in 2009.

# Bias due to missing data

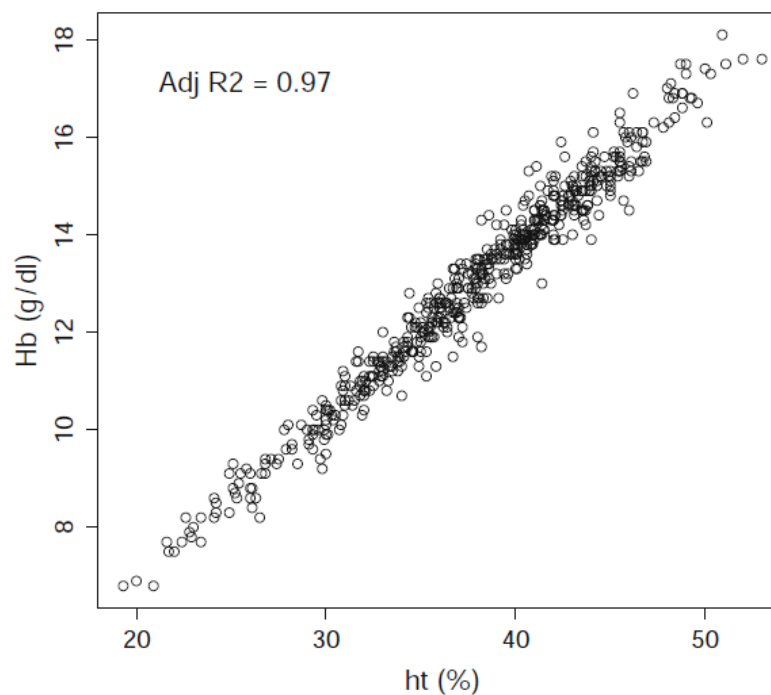


**Fig. 7.1** Effect of missing values on estimated regression coefficients  $\beta_1$  and  $\beta_2$  in the model  $y \sim X_1 + X_2$ . Original data are marked as “dot” and “dash” for  $X_1$  and  $X_2$ , respectively. Complete data under MCAR, MAR, and MNAR are marked with a *circle*. Plots show results for  $n = 500$ ; expected values for  $\beta_1$  and  $\beta_2$  are shown under the graphs (based on  $n = 100,000$ )

Steyerberg. Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating. Springer in 2009.



# Imputation



**Fig. 7.3** Correlation between haematocrit (ht) and haemoglobin (Hb) in 566 patients with traumatic brain injury. The final imputation model included ht ( $p < 0.001$ ) and gender ( $p = 0.01$ ), with  $R^2$  of 0.97

Steyerberg. Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating. Springer in 2009.

# Imputation

- **Sample random normal values**
  - Only external information is used
- **Conditional mean** with a single imputation
  - Predictor data only is used
- **Single imputation** with a random draw from the predictive distribution from a imputation model
  - Predictor data and outcome data are used
- **Multiple imputation** with a random draw from the predictive distribution from an imputation model
  - Predictor data and outcome data are used



# Imputation

- Choosing the imputation
  - Imputation model aims to approximate the true distributional relationship between the unobserved data and the available information
  - Two modelling choices usually have to be made:
    - the form of the model (e.g. linear, logistic, polytomous)
    - and the set of variables that enter the model, including potential transformations of predictors.
  - Truncate imputed values, so that they remain within a plausible range
  - Always include all predictors and the outcome of the final model, consider auxiliary predictors.

# Imputation

- Multiple Imputation
  - In multiple imputation (MI), missing values are imputed  $m$  times using  $m$  independent draws from an imputation model.
  - This means that for each variable with missing data, a conditional distribution for the missing data can be specified given other data
  - $m$  completed data sets are created instead of a single completed data set. Missing values are imputed  $m$  times using  $m$  independent draws from an imputation model.

















# fim

Session 11  
Dealing with missing data

Pedro E A A do Brasil  
[pedro.brasil@ini.fiocruz.br](mailto:pedro.brasil@ini.fiocruz.br)  
2018

