

## 1 DATA S1 MISSING DATA IMPUTATION METHODS

### 2 *Imputation using random forests*

3 The first approach uses random forests to impute missing values and was implemented in the R  
4 package `missForest` which can impute continuous, categorical and mixed data. Random forests do  
5 not have any assumptions of underlying distributions, can accommodate high dimension data and can  
6 handle complex interactions and non-linear relationships among variables (Stekhoven and Bühlmann,  
7 2012). This package was chosen as it does not have any assumptions of the underlying distributions of the  
8 traits, can deal with mixed data and does not require tuning of the parameters (Stekhoven and Bühlmann,  
9 2012). To estimate the unknown variables, `missForest` fits a random forest to the observed part of  
10 each variable with missing data and uses the trained model to predict the missing data. To create a random  
11 forest, many decision trees are trained with each tree only seeing a random subset of the data. Decision  
12 trees classify the data one variable at a time, where the variable becomes the node and the data separated  
13 into groups based on the nodes criteria. To create the final forest, all of the trees are averaged to generate  
14 output probabilities. The random forest was implemented using the R package `randomForest` (Liaw  
15 and Wiener, 2002).

### 16 17 *Imputation using Multiple Correspondence Analysis (MCA)*

18 `missMDA` (Josse and Husson, 2016) uses iterative MCA algorithms, which consists of an initialization  
19 step in which missing values in the indicator matrix are imputed by initial values, such as the proportion of  
20 the category. This initialization for categorical variables is equivalent to mean imputation for continuous  
21 variables. Then, MCA is performed on the imputed indicator matrix to obtain an estimation of parameters,  
22 and missing values are imputed using the fitted values.

### 23 24 *Imputation using Multivariate Imputation by Chained Equations (MICE)*

25 MICE is a popular R data imputation package that offers a range of methods to impute data (van Bu-  
26 uren and Groothuis-Oudshoorn, 2011). For our purposes we used the option from MICE that handles  
27 factor variables with more than 2 categories. MICE imputes data in three steps. First, plausible values  
28 are drawn from a distribution that is specifically modelled for each missing value. This is important as  
29 in other imputation packages and methods it is often assumed that all of the missing values come from  
30 the same distribution. Second,  $m$  imputed datasets are created based on the missing value distributions.  
31 In each of these datasets, the non-missing data remain identical, while imputed values may differ. A  
32 value of  $Q$  is determined for each dataset which allows the uncertainty in each estimation to be calculated.  
33 Finally, estimates of missing data are pooled, and since  $Q$  is approximately normally distributed, the mean  
34 is imputed as the missing value. Since the estimate comes from a normal distribution we are also able  
35 to estimate its variance. For our case, we only have nominal variables, therefore imputing missing data  
36 using a normal distribution is inappropriate. Instead a polytomous regression is used (Brand, 1999). This  
37 allows the underlying relationships between variables to vary non-linearly. This type of regression can  
38 also handle nominal variables with different category levels.

## 40 REFERENCES

- 41 Brand, J. (1999). *Development, implementation and evaluation of multiple imputation strategies for the*  
42 *statistical analysis of incomplete data sets*. PhD thesis, University of Rotterdam.
- 43 Josse, J. and Husson, F. (2016). `missmda`: a package for handling missing values in multivariate data  
44 analysis. *Journal of Statistical Software*, 70(1):1–31.
- 45 Liaw, A. and Wiener, M. (2002). Classification and regression by randomforest. *R news*, 2(3):18–22.
- 46 Stekhoven, D. J. and Bühlmann, P. (2012). `Missforest`: non-parametric missing value imputation for  
47 mixed-type data. *Bioinformatics*, 28(1):112–118.
- 48 van Buuren, S. and Groothuis-Oudshoorn, K. (2011). `mice`: Multivariate imputation by chained equations  
49 in r. *Journal of statistical software*, 45(3):1–68.