



Swiss Federal Institute of Technology Zurich

Seminar for  
Statistics

Department of Mathematics

---

Semester Paper

Fall 2015

---

David Pham

# Missing Data: Empirical Comparison between Imputation and Nearest Neighbors Algorithms

---

Submission Date: February 15th 2016

---

Adviser: Dr. Martin Maechler



To the *R* community and *ESS* developers for their contribution.



# Abstract

Incomplete or missing data are common in scientific works, and the most common solution to cope with them is the simply to discard them. However, this might lead to bias in the conclusion. This semester paper summaries modern methods and offers an empirical comparison of the packages *amelia*, *imputeKnn*, *mi*, *mice*, *softimpute* with the statistical software R.

## Contents

<b>1</b>	<b>Theoretical Background</b>	<b>1</b>
1.1	Mechanism of missingness . . . . .	1
1.2	Statistical completion . . . . .	2
1.3	Algorithmic completion . . . . .	3
<b>2</b>	<b>Empirical Comparison of Imputation Methods</b>	<b>7</b>
2.1	Data set and R packages . . . . .	7
2.2	Methodology . . . . .	7
2.2.1	Simulation of missingness . . . . .	8
2.2.2	Ranking methods . . . . .	9
2.3	Implementation constraints . . . . .	9
2.4	Results . . . . .	10
2.5	Open questions . . . . .	10
<b>3</b>	<b>Conclusion</b>	<b>17</b>
	<b>Bibliography</b>	<b>19</b>

## List of Figures

2.1	Relative ranking of imputation quality of the tuning parameters of soft-Impute and impute.knn. For impute.knn, the number of neighbors is the tuning parameters, whereas for softimpute, it is the maximum rank and estimation method of the output matrix. . . . .	11
2.2	Rankings of imputation methods on the FLAS data set grouped by missing rate, under the MCAR mechanism with missing rate. Labels in the boxes provide the missing rate. . . . .	13
2.3	SMSE for selected missingness rate with MCAR against imputation methods. . . . .	14
2.4	SMSE of imputation methods on the FLAS data set grouped by missing rate, under the MAR mechanism. Labels in the boxes provide the missing rate. . . . .	15

## List of Tables

2.1	FLAS data set, summary of numerical variables . . . . .	8
2.2	FLAS data set, summary of factor variables . . . . .	8





# Chapter 1

## Theoretical Background

This chapter provides an overview and an intuition on the field of missing data. It follows mainly [Schafer and Graham \(2002\)](#), [Little and Rubin \(2002\)](#), [Van Buuren \(2012\)](#), with some impute from [Wikipedia \(2015\)](#), [Matloff \(2015\)](#), [Gelman and Hill \(2006\)](#), [Troyanskaya, Cantor, Sherlock, Brown, Hastie, Tibshirani, Botstein, and Altman \(2001\)](#). This chapter begins with a short description on the nature missingness, then describes several procedures in order to handle missing data.

### 1.1 Mechanism of missingness

[Van Buuren \(2012\)](#) describes two concepts helping us to understand how to solve the problem of missing data: intentional and unintentional missingness, as well as unit and item missingness. The experimenter can decide to not measure all possible variable in an experiment and encode his decisions as missing observations. This is a reasonable decision if the cost of measuring variables is material and unnecessary for some experimental case, such as in medical experimentation. However, it might also happen that the experimenter could not measure some variable, e.g. when a respondent to a survey refuse to answer to some question. In this case, the missingness is named unintentional. The second concept is often missingness is about unit and items: one says a unit is missing when none of the variables of interest could be measured, whereas item refers to some variable missing.

In order to complete missing data, assumptions need to be taken about the underlying mechanism creating missing observations: missing completely at random (MCAR), missing at random (MAR) and missing not at random (MNAR).

**Notation** Let  $Y \in \mathbb{R}^{n \times p}$  be the data matrix containing missing data for  $n$  observations with  $p$  variables,  $R = (r_{ij})_{i=1, j=1}^{n, p} \in \{0, 1\}^{n \times p}$  denotes the response  $y_{ij}$  (i.e.  $r_{ij} = 1$  is  $y_{ij}$  is observed, and is 0 otherwise).  $Y_{obs}$  and  $Y_{mis}$  denote the observation that is observed, respectively, missing, such that  $Y = (Y_{obs}, Y_{mis})$ . Note that we always observed  $R$  whereas we usually do not have  $Y_{mis}$ .

**MCAR** The data are said to be *MCAR* if

$$P(R = 0 \mid Y_{obs}, Y_{mis}) = P(R = 0),$$

or equivalently

$$P(Y = y \mid R = i) = P(Y = y), \quad i \in \{0, 1\}, \quad y \in \mathbb{R}^{n \times p}.$$

It means the probability of being missing depends does not depends on the actual value of  $Y$ .

**MAR** For multiple imputation, one requires only  $R \perp Y_{mis}$ , that is

$$P(R = 0 \mid Y_{obs}, Y_{mis}) = P(R = 0 \mid Y_{obs}),$$

that is other observed variables impact of the probability of missingness but the missing mechanism only depends on the observed variables and not the actual missing value. In this case, we say the data  $Y$  are *MAR*.

**MNAR** The data are MNAR if

$$P(R = 0 \mid Y_{obs}, Y_{mis})$$

can not be simplified. It essentially means that the rate of response depends on the actual value of the missing observations. The standard example is the survey about salary where people with extremely high salary tends to hide their earnings.

Modern statistical technique can handle MNAR and MAR cases, whereas simple technique only MCAR, which is quite restrictive.

## 1.2 Statistical completion

**Complete case analysis** Unfortunately, One of the most used technique to cope with missing data: the researcher only keeps observation that are complete. This might lead to valid analysis, as the method does not introduce any bias if the missing values are uniformly distributed. Nevertheless, this methodology can not work in modern settings where the probability of one missing variable is quite high: Many data points would be discarded.

**Pairwise deletion** This methods improve from the previous one by deleting observations only if the variable which is missing must be used in the model. This is typically relevant for computing correlation for example, although some care must be taken in this case.

**Single imputation** The data matrix is sorted according to some order, *last observation carried forward* is the method of replacing the missing value with last valid value. The missing value can also be replaced with the mean of the other observations, however, correlations are attenuated. Regression imputation use the other variables as predictors to replace the missing value, although precision is misleadingly augmented, hence does not reflect the statistical errors of the missing data. This problem is partially solved by multiple imputation.

**Multiple imputation** Under the MAR assumption, the multiple imputation (MI) is similar to bootstrapping method: the distribution of each variable conditional and the others is fitted, then in case of missing value, a sample is drawn from this distribution. The desired statistics are averaged except for the standard error which is constructed by adding the variance of the imputed data and the within variance of each data set. The last step solves the problem of understating uncertainty. Standard errors reflect missing-data uncertainty and finite-sample variation.

More precisely, in the one-dimensional case, if the sample is large enough so that the estimator  $Q$  follows a Gaussian distribution, then the estimate  $\hat{Q}$  and the standard error  $T$  can be computed from the estimates of  $(Q^j, U^j)_{j=1}^m$ ,  $Q^j$ , respectively,  $U^j$  being the fitted value of  $Q$ , respectively the standard error, for data sets  $j$ :

$$\begin{aligned}\hat{Q} &= m^{-1} \sum_{j=1}^m Q^j, \\ \hat{U} &= m^{-1} \sum_{j=1}^m U^j, \\ B &= (m-1)^{-1} \sum_{j=1}^m (Q^j - \hat{Q})^2, \\ T &= \hat{U} + (1 + m^{-1})B.\end{aligned}$$

For confidence interval, the Student's  $t$  approximation can be used with the degree of freedom given by

$$\nu = (m-1) \left[ 1 + \frac{\hat{U}}{(1 + m^{-1})B} \right]^2.$$

The estimated rate of missing information for  $Q$  is approximately  $\tau/(\tau+1)$  where  $\tau = (1 + m^{-1})B/\hat{U}$ , the relative increase in variance due to non-response. See [Schafer \(1997\)](#) for more cases.

An advantage of MI is the number of need imputation: the efficiency based on  $m$  samples relative o an infinite number is  $(1 + \lambda/m)^{-1}$ , where  $\lambda$  is the rate of missing information, which measures the increase in the large-sample variance of a parameter estimate due to missing values.  $m = 20$  is often good in practice.

Obviously, the missing values problem is dealt before the analysis with MI, in contrast with maximum likelihood estimation. The danger from MI is the ability to use different models for imputation and analysis, which might lead to inconsistency.

### 1.3 Algorithmic completion

For this particular section, the data matrix will be denoted as  $X$ , in order to remain consistent with the literature.

**Singular value decomposition** Singular values of a matrix  $X$  are the square root of the non-negative eigenvalues of  $X^T X$ . Singular value decomposition (SVD) is provided by

$$\hat{X}_J^c = U_J D_J V_J^T, \quad (1.1)$$

where  $D_J \in \mathbb{R}^{N \times p}$  is a diagonal matrix containing the leading  $J < p$  singular values of  $X^c$  and  $V_J \in \mathbb{R}^{p \times p}$  and  $U_J \in \mathbb{R}^{N \times N}$ , the corresponding orthogonal matrix of  $J$  right and left singular vectors.  $\hat{X}^c$  is the nearest matrix of  $X^c$  among matrices with rank  $J$  with respect to the sum of squares norm  $\|A\|^2 = \text{tr}(AA^T)$ .

If  $x_i$  is any row of  $X^c$ , consider the regression of the  $p$  values in  $x_i = (x_{i1}, \dots, x_{ip})^T$  on the eigen-vectors  $v_1, \dots, v_J$ , each  $p$  dimensional vectors. The regression solves

$$\min_{\beta} \|x_i - V_J \beta\|^2 = \min_{\beta} \sum_{l=1}^p (x_{il} - \sum_{j=1}^J v_{lj} \beta_j)^2,$$

with solution  $\hat{\beta} = (V_J^T V_J)^{-1} V_J^T x = V_J^T x$  (since  $V_J$  is orthogonal) and orthogonal values  $\hat{x}_l = V_l \hat{\beta}, l \in \{1, \dots, J\}$ . Thus, according to Equation (1.1),  $X^c V_J = U_J D_J$  gives all the (transposed) regression coefficients for all the rows and  $\hat{X}^c = U_J D_J V_J^T$  all the fitted values. Hence, once the matrix  $V_J$  is computed, SVD approximate each row of  $X^c$  by its fitted vector obtained by regression (or projection) on  $V_J$ . This suggest for a row  $x_i$  of  $X^m$  with some missing components, they could possibly be imputed from

$$\min_{\beta} \sum_{l=1}^p 1(R_{il} = 1) (x_{il} - \sum_{j=1}^J v_{lj} \beta_j)^2,$$

where  $R_{il}$  is the response indicator of  $x_{il}$ .

The imputation procedure is described as the following.

- i.) Compute the SVD of  $X^c$  and keep  $V_J$ .
- ii.) For a row  $x^*$  with missing element, compute

$$\hat{\beta}^* = (V_J^{*T} V_J^*)^{-1} V_J^{*T} x^*,$$

where  $V_J^*$  is the shortened version of  $V_J$  with the appropriate rows removed (corresponding the missing elements of  $x^*$ ). Note  $V_J^*$  no longer has orthogonal columns.

- iii.) The predictions of the missing elements are  $V_J^{(*)} \hat{\beta}^*$  where  $V_J^{(*)}$  is the complement in  $V_J$  of  $V_J^*$ .

Usually, the data matrix is centered before SVD, however, for missing data, an intercept has to be fitted and a method based simulation is provided afterwards. The previous methods usually discards a great number of data, particularly when  $p \gg N$ . In contrast, the next iterative procedure circumvent the problem at the cost of more computation.

- (1) Set  $X^*$  as  $X$  with all missing values filled by the mean of their row.
- (2) Solve the problem

$$\min_{V_J, D_J, U_J} \|X^* - m 1^T - U_J D_J V_J^T\|^* \quad (1.2)$$

where  $\|\cdot\|^*$  is the sum of squares of all non-missing elements and  $m \in \mathbb{R}^N$  is the row means of  $X^*$ .

- (3) Predict the missing values of  $X$  with the fitted values.
- (4) Reset  $X^*$  as  $X$  with the missing values replaced by the result of previous step.
- (5) Repeat steps 2-5, until the size of the relative update of the missing values become negligible.

According to [Hastie, Tibshirani, Sherlock, Eisen, Brown, and Botstein \(1999\)](#), only 6 iterations are necessary. Interestingly, the solution of Equation (1.2) is a fixed point, i.e. if missing values are filled, and the SVD algorithm is executed on the complete matrix, the solution remains identical.

**K-nearest neighbors completion** [Troyanskaya et al. \(2001\)](#) presents the other end of the spectrum in term of data usage: *K-nearest neighbor averaging*. The algorithm is described as following.

- i.) Computed the Euclidean distance between  $x^*$  and all the rows in  $X^c$ , using only those co-ordinates not missing in  $x^*$ . Identify the  $K$  closest observations.
- ii.) Impute the missing coordinates of  $x^*$  by averaging the corresponding coordinates of the  $K$  closest with weights proportional to their distances to  $x^*$ .

Empirically, the number of neighbors  $K$  between 5 to 10 is often a good choice for most data set.



## Chapter 2

# Empirical Comparison of Imputation Methods

### 2.1 Data set and R packages

**Data set** The FLAS data set is studied in [Schafer \(1997\)](#) and is a great candidate for the simulation study. The data were collected in 1987 to investigate the impact of the Foreign Language Attitude Scale (FLAS), a new measure, for predicting success in learning new foreign languages. Tables [2.1](#) and [2.2](#) provide a short summary of data. The MICE and `mi` packages have been applied to the original data set to create an artificial complete one. The latter is used as baseline to compare the imputations by different methods.

**R Packages** Five R packages were selected in the study.

**Amelia** implemented by [Honaker, King, and Blackwell \(2011\)](#) provides bootstrapping methods and EM algorithm for multiple impute analysis.

**impute** authored by [Hastie, Tibshirani, Narasimhan, and Chu \(1999\)](#) uses the `impute.knn` function to provide nearest neighbors imputation.

**mi** created by [Gelman and Hill \(2011\)](#) implements the multiple imputations in a Bayesian framework.

**mice** written by [van Buuren and Groothuis-Oudshoorn \(2011\)](#) allows the user to impute values with chained equations.

**softImpute** distributed by [Hastie and Mazumder \(2015\)](#) uses singular value decomposition (or a version thereof) to complete data sets.

Each of these packages offers a function for completing data set. The `simslapar` packages from [Hofert and Maechler \(2015\)](#) provides a stable framework to conduct the simulation and gather the output from parallel simulations.

### 2.2 Methodology

Let  $\mathcal{M}$  denote the set of imputation methods and let  $Y \in \mathbb{R}^{n \times p}$  be a complete data matrix.

Table 2.1: FLAS data set, summary of numerical variables

Statistic	N	Mean	St. Dev.	Min	Max
FLAS	279	82.487	14.026	28	110
MLAT	230	24.257	6.256	9	40
vSAT	245	501.514	91.162	210	790
mSAT	245	564.249	88.707	320	800
eng	242	53.950	15.402	19	113
HGPA	278	2.750	0.617	0.500	3.990
CGPA	245	3.294	0.477	2.000	4.000

Table 2.2: FLAS data set, summary of factor variables

Statistic	N	Factors				
Age	268	-19	20+			
		124	144			
Sex	278	M	F			
		152	126			
Number of prior foreign language	268	none	1-2	3+		
		71	73	124		
Prior Language	279	french	spanish	german	russian	
		67	78	114	20	
Grades	232	F	D	C	B	A
		1	5	22	79	125

The experience requires to chose a complete data matrix, then to replace some of its entries with missing values and then to rank the imputation methods. In order to cope with the randomness,  $N_{sim}$  simulations are performed and then aggregated.

### 2.2.1 Simulation of missingness

Two types of missing were implemented: MCAR and MAR. The former is quit straightforward to implement: For a given missingness rate  $p \in (0, 1)$ , the response<sup>1</sup>  $r_{ij}$  of the value  $Y_{ij}$  follows a binomial distribution with probability  $p$ . Implementation of MAR usually make assumptions on the underlying multivariate distribution is a little more involved. Nevertheless, a mechanism based of the empirical distribution of the missingness pattern can also been used.

The original data set contained pattern of missingness and for a missing rate  $q$ , one could sample the patterns from it with the multinomial distribution until the desired rate  $q$  is reached. The patterns are then randomly assigned to our completed data set. This

<sup>1</sup>Recall that  $r_{ij}$  is 0 if  $y_{ij}$  is missing and 1 otherwise.



method has the disadvantage that it can not attain any missingness rate  $q \in (0, 1)$ , as one can only assign one pattern per observation.

For our study, it implies that for the MAR method, only data set with a missingness rate below 30% were analyzed.

### 2.2.2 Ranking methods

For a simulated data matrix  $Y^l = (y_{ij}^l)_{i,j=1}^{n,p}$ ,  $l \in \{1, \dots, N_{sim}\}$  with missing values, its associated response matrix  $R_{ij}^l = (r_{ij}^l)_{i,j=1}^{n,p}$  and an imputation method  $m \in \mathcal{M}$  with predicted values  $\hat{y}_{ij}^l$  if  $y_{ij}^l$  is missing. For *numerical* variables, the scaled mean squared error (SMSE),

$$\text{SMSE}_{l,j}^m = \left\{ \sum_{i=1}^n 1(r_{ij}^l = 0) \right\}^{-1} \sum_{i=1}^n \left( \frac{\hat{y}_{ij}^l - y_{ij}^l}{\mu_j} \right)^2 \cdot 1(r_{ij}^l = 0), \quad j \in \{1, \dots, p\}, \quad m \in \mathcal{M}, \quad (2.1)$$

where  $\mu_j = n^{-1} \sum_{i=1}^n y_{ij}$ , is used to assess the quality of imputation methods. For *factors* variables, the conservative 0 – 1 loss is employed. To aggregate it across the columns of the data matrix, for each column, the methods are ranked by SMSE, then these ranks are summed by imputation method. More precisely, the score  $s_l^m$  of the imputation methods  $m \in \mathcal{M}$  can be expressed as

$$s_l^m = \sum_{j=1}^p \sum_{\nu \in \mathcal{M}} 1(\text{SMSE}_{l,j}^m \leq \text{SMSE}_{l,j}^\nu). \quad (2.2)$$

The scores  $s_l^m$  are then used to assess the performance of the imputation methods.

## 2.3 Implementation constraints

**Data type for soft impute and nearest neighbors** The implementation of two methods did not allow for factors variable. Although it is not a paramount task to transform the factor into numerical data type, some care should be taken to make conversion correctly.

**Nearest neighbors imputation** It appears that the `impute.knn` method from the `impute` package from the *bioconductor* repository causes *segmentation fault* (with the underlying FORTRAN code) when the number of neighbors is either too high with respect to the available data.

**Collinear dimensions** If the data matrix  $Y$  has collinear variables, then some challenges might occur when variables are collinear. More precisely, although multiple imputation techniques try to estimate

$$Y_j \mid Y_{k_1}, \dots, Y_{k_j},$$

using regression models, their results might be unstable if  $Y_{k_i}$ ,  $i \in \{1, \dots, k_j\}$  are linearly dependent. Said differently, as regression parameters depends on the quality of the inversion the data matrix, but a matrix with collinear columns has an unstable inverse<sup>2</sup>. The **Amelia** package is the most prone to this issue, although **MICE** and **mi** algorithms fail to converge at some point when  $Y_{k_i}$  exhibit high collinearity.

**Timing** Usually, CPU time, i.e. time spent by the processor on the *R* process, is measured to evaluate the speed performance. Nonetheless, a trend has emerged for packages implement pretty good parallelism, i.e. packages implement the parallelism procedures themselves. It leads to underestimated human elapsed time as most of the computational burden is performed by sub-process.

**Tuning parameters** For the packages **softImpute** and **impute**, some defaults parameters for the imputation methods are provided. Figure 2.1 displays how the quality of the imputation evolves with the parameters. The **impute.knn** function, the quality of the inference grows with the number of neighbors. The **softImpute** function unexpectedly offers good default parameters, even if some restrain should be kept when the missing rate is high.

## 2.4 Results

Figure 2.2 summarize the results under the MCAR missing mechanism. **softImpute** and **MICE** were the only package able to cope with a quite high missing rate ( $p \geq 0.7$ ). **impute.knn** from the **impute** package, when it works, is almost always the method with the best score, as defined in Equation (2.2). **MICE** offers a good balance between speed, robustness and quality of imputations, but the methods depends on the linear dependency of the columns of the data matrix. Although **softImpute** can cope with almost any type of data matrix, its inferences are sub-par with the other methods. Some further analysis might be needed to confirm this results. **Amelia** is the package which is the least able to cope with a random impute matrix: routines failed without exception with any missing rate above  $p > 0.3$ . The **Amelia** and **mi** packages use by default parallel back-ends to perform their computations. However, they are the slowest methods in terms of elapsed time. This might be overcome by setting the number of iteration to a lower threshold. Nonetheless, this is not recommended as convergence is not guaranteed when data input might have dimensions with strong linear dependency.

Figure 2.3 shows the measure from Equation (2.1) for missingness rates  $p \in \{0.2, 0.5, 0.7\}$ . For numerical data, all methods but **softImpute** have the same order of errors with K-nearest neighbors being slightly better than the others. However, the latter method does not allow for inference of the imputation error. As Figure 2.4 shows, the previous statements are not impacted by changing the missingness mechanism to (our implementation of) MAR.

---

<sup>2</sup>The matrix is so-called *ill conditioned*.

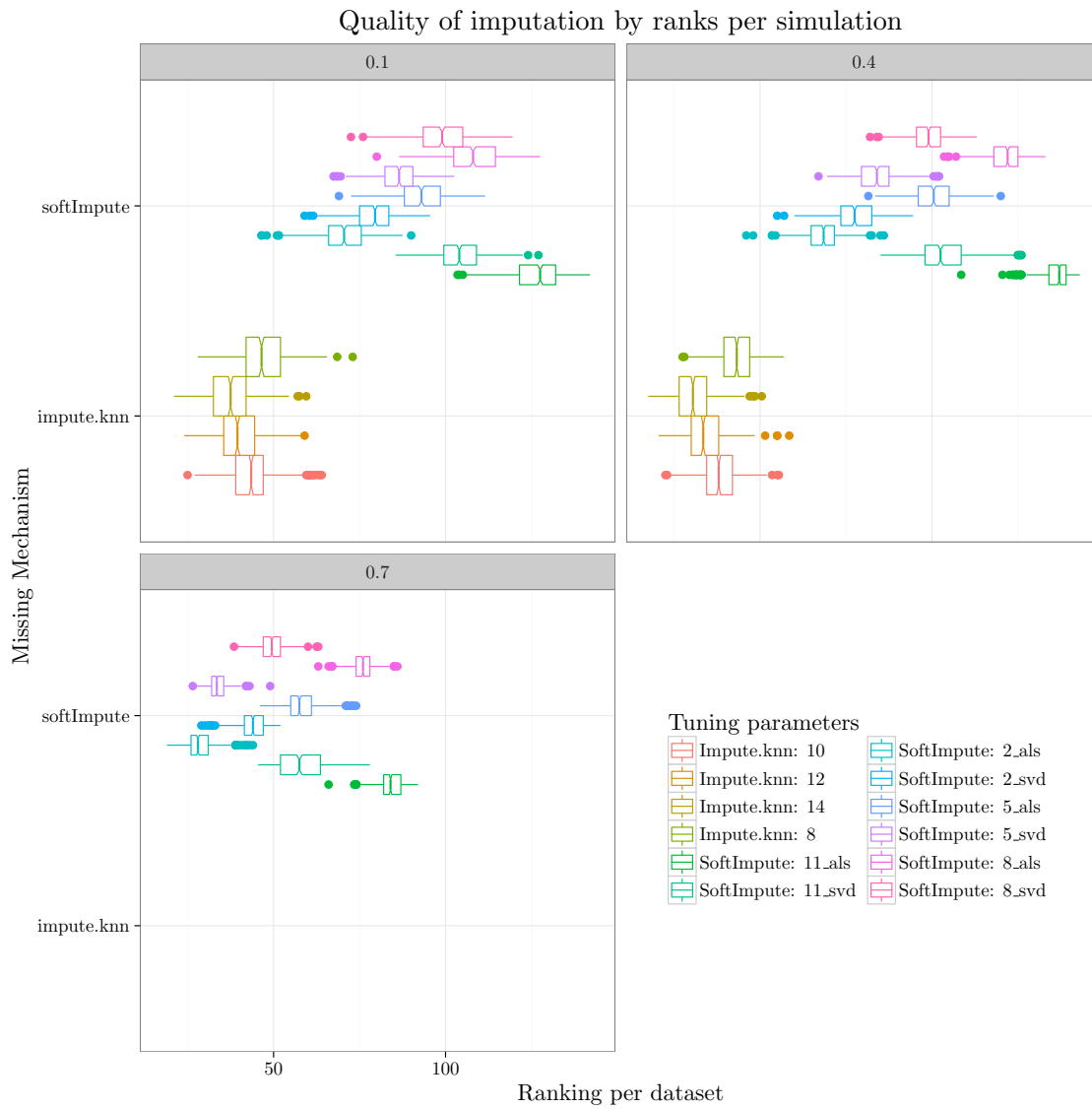


Figure 2.1: Relative ranking of imputation quality of the tuning parameters of softImpute and impute.knn. For impute.knn, the number of neighbors is the tuning parameters, whereas for softimpute, it is the maximum rank and estimation method of the output matrix.

## 2.5 Open questions

Heuristically, the quality of models output normally depends on the amount of available data. In the missing data framework, precision are needed for this notion: Is it the number of complete observations, the number of non-missing values, or a mixture of both? Interactions between the number of observation  $n$ , the number of dimensions  $p$  and imputation methods with their optimal parameters are left unanswered with this work. In order to answer this question, a multivariate sampling mechanism should be devised and tested.

Moreover, are there any reasonable solutions which can be applied to overcome collinearity? One could cluster the similar dimension and then pick one randomly to create an imputed value. Nevertheless, such solutions were not yet implemented.

Additionally, this small simulation study has been applied to certain data set and it should be interesting to repeat the experience with other real data.

Finally, the concern of this work has been to apply imputation methods to retrieve potential candidate values for inference, it would have been interesting to verify the quality of inferences performed with the imputed data set, for example multiple imputation against nearest neighbors.

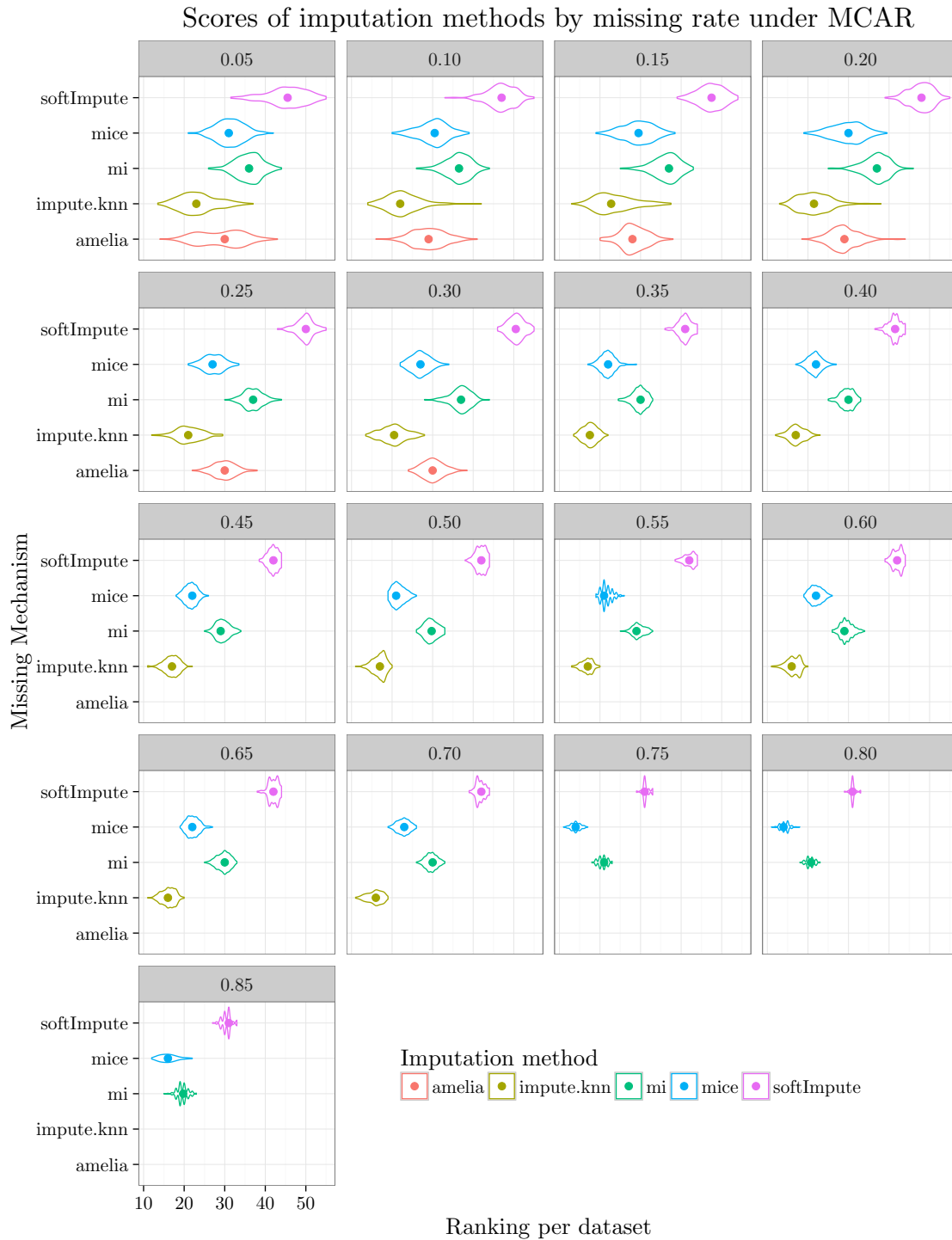


Figure 2.2: Rankings of imputation methods on the FLAS data set grouped by missing rate, under the MCAR mechanism with missing rate. Labels in the boxes provide the missing rate.

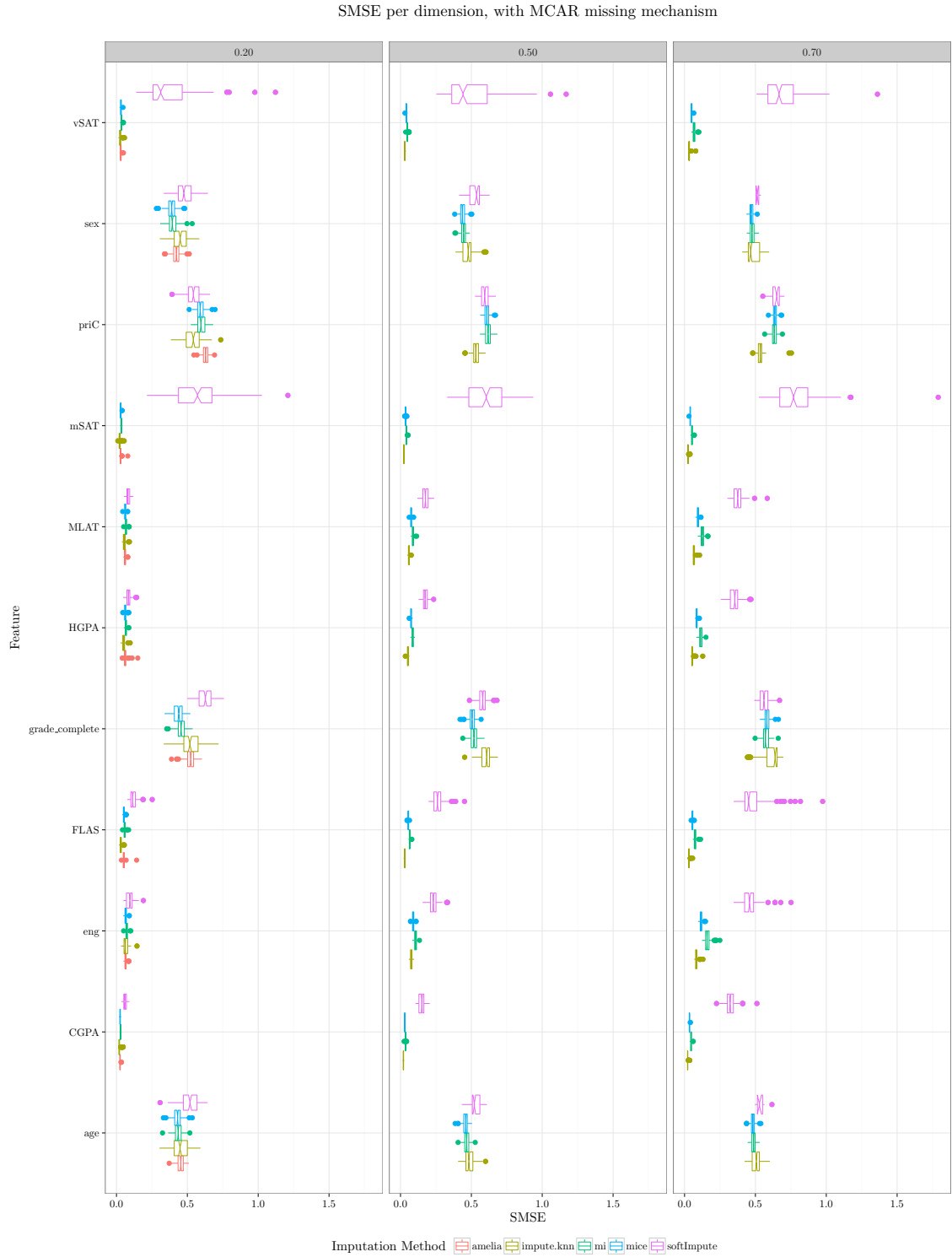


Figure 2.3: SMSE for selected missingness rate with MCAR against imputation methods.



Figure 2.4: SMSE of imputation methods on the FLAS data set grouped by missing rate, under the MAR mechanism. Labels in the boxes provide the missing rate.





## Chapter 3

# Conclusion

In this semester papers, part of theory of data completion is reviewed. The main work is devoted to build a framework to test several R packages for imputation methods on the FLAS data set. With this data set, it appears that K-nearest neighbor imputation works fairly well, although its implementation might throw frustrating low level errors (segmentation faults). Except for `softImpute`, all methods have the same order of error (distance between the imputed and true value). In practice, they could unfortunately not be use as black-box as most data matrix have colinear dimensions which constitutes an issue for all the algorithms based on regressions.

Finally, as departing words, one should not forget why these technique exists. From [Schafer and Graham \(2002\)](#),

With or without missing data, the goal of a statistical procedure should be to make valid and efficient inferences about a population of interest – not to estimate, predict, or recover missing observations nor to obtain the same results that we would have seen with complete data.



# Bibliography

- Gelman, A. and J. Hill (2006). *Data analysis using regression and multilevel/hierarchical models*. Cambridge University Press.
- Gelman, A. and J. Hill (2011). Opening windows to the black box. *Journal of Statistical Software* 40.
- Hastie, T. and R. Mazumder (2015). *softImpute: Matrix Completion via Iterative Soft-Thresholded SVD*. R package version 1.4.
- Hastie, T., R. Tibshirani, B. Narasimhan, and G. Chu (1999). *impute: impute: Imputation for microarray data*. R package version 1.42.0.
- Hastie, T., R. Tibshirani, G. Sherlock, M. Eisen, P. Brown, and D. Botstein (1999). Imputing missing data for gene expression arrays.
- Hofert, M. and M. Maechler (2015). *simsalapar: Tools for Simulation Studies in Parallel with R*. R package version 1.0-5.
- Honaker, J., G. King, and M. Blackwell (2011). Amelia II: A program for missing data. *Journal of Statistical Software* 45(7), 1–47.
- Little, R. and D. Rubin (2002). Statistical analysis with missing data.
- Matloff, N. (2015). New r software/methodology for handling missing data.
- Schafer, J. L. (1997). *Analysis of incomplete multivariate data*. CRC press.
- Schafer, J. L. and J. W. Graham (2002). Missing data: our view of the state of the art. *Psychological methods* 7(2), 147.
- Troyanskaya, O., M. Cantor, G. Sherlock, P. Brown, T. Hastie, R. Tibshirani, D. Botstein, and R. B. Altman (2001). Missing value estimation methods for dna microarrays. *Bioinformatics* 17(6), 520–525.
- Van Buuren, S. (2012). *Flexible imputation of missing data*. CRC press.
- van Buuren, S. and K. Groothuis-Oudshoorn (2011). mice: Multivariate imputation by chained equations in r. *Journal of Statistical Software* 45(3), 1–67.
- Wikipedia (2015). Imputation (statistics).



# Declaration of Originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor .

**Title of work** (in block letters):

**Authored by** (in block letters):

*For papers written by groups the names of all authors are required.*

**Name(s):**

**First name(s):**

Muster	Student

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the Citation etiquette information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work .
- I am aware that the work may be screened electronically for plagiarism.
- I have understood and followed the guidelines in the document *Scientific Works in Mathematics*.

**Place, date:**

**Signature(s):**

Zurich August 19th 2009	bla

*For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.*