

Exact solution of the donor-limited nearest neighbor hot deck imputation problem

Jan Pablo Burgard*, Sven de Vries†, Ulf Friedrich ‡, and Dennis Kreber †§

Abstract. Data quality in population surveys suffers from missing responses. We use combinatorial optimization to create a complete and coherent data set. The methods are based on the widespread nearest neighbor hot deck imputation method that replaces the missing values with observed values from a close unit, the so-called donor. As a repeated use of donors may lead to distortions in the statistics, an additional constraint on the maximum number of donor re-uses is introduced. It is shown how this problem can be modeled as a maximum weighted b -matching problem. Based on this theoretical analysis, we implement a cost scaling, network simplex and successive shortest path algorithm for computing a globally optimal solution. These outperform the available and presently used implementations in terms of solution quality and running time. Statistical imputation is presented as a novel application of Operations Research, the available methodology is improved by a concise theoretical treatment, and fast implementations are provided.

Key words. Survey Data, Combinatorial Optimization, b -Matching, Missing Data, Nearest Neighbor Hot Deck Imputation

AMS subject classifications. 62D05, 62-07, 90C27, 90C35, 90C90

1. Introduction. The problem of missing data is omnipresent in survey research and occurs in the two main forms of *unit non-response*, when the information for a whole unit is missing, and of *item non-response*, when single values for some units are missing [30]. Missing data cause severe problems in the statistical evaluations such as biased estimates. Moreover, incomplete data sets inhibit the direct application of estimation methods that require complete data sets, see for instance [32]. In survey statistics, it is therefore a central task to replace the missing data with appropriate values. This replacement process is called *imputation*.

In the unit non-response case, weighting adjustments are chosen to account for a possible non-response bias [13]. In the item non-response case often *single imputation* methods are used and are even more often preferred in official statistics, see [12]. In official statistics, generally a complete data set is desired for computing the necessary population figures. This is due

*Economic and Social Statistics, Trier University, 54286 Trier, Germany (burgardj@uni-trier.de)

†Mathematics, Trier University, 54286 Trier, Germany (devries@uni-trier.de)

‡Chair of Operations Research, Technical University of Munich, Arcisstr. 21, 80333 Munich, Germany (ulf.friedrich@tum.de)

§Germany, DFG-RTG Algorithmic Optimization, Trier University, 54286 Trier, Germany (kreberd@uni-trier.de)

30 to the fact that a complete data set automatically yields hierarchical consistent population
31 figures, e.g., the sum of counties-figures equals the figure on the whole population.

32 Only one value is imputed in each missing field in single imputation. Hot deck single
33 imputation [12] is a imputation method that maintains plausibility of records to a greater
34 extend. The *hot deck* consists of units that are chosen to be possible donors, e.g., all units
35 without missing values. An extensive overview of several imputation methods and how they
36 compare to the hot deck imputation is given in [22], [2], and [26]. We focus on the *nearest*
37 *neighbor hot deck imputation*, which chooses the donor for a unit with missing value(s) in
38 terms of minimizing a predefined distance [21]. The method plays an important role in official
39 statistics and was used in censuses conducted in Canada, the United Kingdom, Italy, Brazil,
40 the USA, New Zealand, Australia and Switzerland [3] and recently in Germany [11].

41 In practice it may be desirable to limit the amount of times a donor may be re-used in a
42 hot deck. In [19] the effect of limiting the number of times a donor may be re-used is studied
43 and concluded that in general a reduced variance can be expected when taking each donor
44 the least possible time. The authors of [18] combine the nearest neighbor aspect with a donor
45 limitation because this imputation method is especially prone to an over-usage of donors.
46 They examine the effects and find that a donor-limited nearest neighbor hot deck imputation
47 is especially advantageous on data sets with a high percentage of missing values. They advise
48 to minimize the sum of *all* distances between donors and recipients. Following up on this,
49 [16], [17] state that this minimization problem can be formulated as a transportation problem.
50 The author proposes two heuristic algorithms that, however, do not guarantee the quality of
51 the computed solution.

52 In contrast, we find a globally optimal solution to the donor-limited problem with an
53 arbitrary donor-specific integer-limit. Based on a graph theoretic model for the problem, we
54 develop exact combinatorial algorithms and compare their speed and stability with respect to
55 the total number of units in the data set.

56 **Contribution and Structure.** The nearest neighbor hot deck imputation is of great signifi-
57 cance in official statistics and an additional donor-limit introduces a beneficial extension to the
58 imputation methodology. Yet, to our knowledge, work on this topic has only been conducted
59 by [16], [17] and [18]. We propose mathematical precise methods for finding a donor-limited
60 nearest neighbor hot deck imputation. Our contributions include the development of models
61 on which we can apply combinatorial algorithms to find globally optimal solutions to the
62 donor-limited hot deck imputation problem. In contrast to the work of [16], [17], we propose
63 using *exact* methods instead of applying heuristics.

64 The presented model gives a more flexible control over possible imputations. Donor limits
65 can be chosen independently for each donor unit and specific imputation pairs can be ex- or
66 included from the model, which allows more control over desired statistical properties.

67 We show that even though we calculate the global optimum, computation times stay
68 at a minimum. In fact, our implementations are considerably faster than the R-Package
69 *HotDeckImputation* [18]. Furthermore, we find that the successive shortest path algorithm
70 often benefits from the structure of the problem graph.

71 In Section 2, a concise mathematical formulation of the nearest neighbor hot deck imputa-
72 tion is presented. It includes the definitions of missing patterns, distances and relations,
73 and the formulation of the mathematical problem. Then, the optimal solution to the problem
74 is derived in Section 3. In the last section, numerical results based on the publicly available
75 data set AMELIA [24] are presented. AMELIA is a synthetic and realistic household data
76 set. The paper concludes with a summary and outlook.

77 **2. Nearest Neighbor Hot Deck Imputation.** Before setting up the optimization prob-
78 lem, we define the donors, receivers and the distances between them in a formal way. Let
79 X_1, \dots, X_Q be subsets of \mathbb{R} , the so-called *attributes*. For $N \in \mathbb{N}$ a *population* $\Omega \in \prod_{j=1}^Q X_j^N$
80 of *size* N is defined as an $N \times Q$ matrix whose values in each column j are elements of the
81 attribute X_j . Even though we allow categorical and dichotomous variables (see A), we assume
82 without loss of generality that they are encoded as values in \mathbb{R} .

83 **Definition 2.1.** Let $n \leq N$ and $q \leq Q$. A submatrix $S \in \mathbb{R}^{n \times q}$ of Ω is called a *complete*
84 *sample*. A *sample* is a matrix $\mathcal{S} \in (\mathbb{R} \cup \{\text{NA}\})^{n \times q}$ in which some elements of S are replaced
85 with the value *NA* for “not available”.

86 Each row (u_{i1}, \dots, u_{iq}) of a sample \mathcal{S} is called a *unit*. We also write $(u_{i1}, \dots, u_{iq}) \in \mathcal{S}$
87 even though \mathcal{S} is not a set in the technical sense. We also omit the row index and just write
88 (u_1, \dots, u_q) or u to denote a unit. A unit u is called a *respondent* if $u_j \neq \text{NA}$ for every j .
89 Otherwise, u is called a *non-respondent*.

90 **Definition 2.2.** Let \mathcal{S} be a sample. A relation \prec on $\mathcal{S} \times \mathcal{S}$ is called a *receiver-donor*
91 *relation* if $(u_1, \dots, u_q) \prec (v_1, \dots, v_q)$ implies that u is a non-respondent and that $v_i \neq \text{NA}$ for
92 all i with $u_i = \text{NA}$. The set $A_{\prec}(\mathcal{S}) := \{(v, u) \in \mathcal{S}^2 : u \prec v\}$ describes all pairs of receivers u
93 and donors v . The donor set D_{\prec} and receiver set R_{\prec} are defined as follows:

$$94 \quad (2.1) \quad R_{\prec} := \{u \in \mathcal{S} : u \prec v \text{ for some } v \in \mathcal{S}\},$$

$$95 \quad (2.2) \quad D_{\prec} := \{v \in \mathcal{S} : u \prec v \text{ for some } u \in \mathcal{S}\}.$$

97 In unambiguous cases the index \prec will be omitted in the above definitions and we write
98 $A(\mathcal{S})$, R , and D instead. The reasoning behind the concept of a receiver-donor relation is of
99 course that $v \in \mathcal{S}$ is a potential donor for $u \in \mathcal{S}$ whenever $u \prec v$ holds. Generally, there
100 is more than one possible donor for a given receiver and one has to choose the in a certain
101 sense *best* donor while respecting possible constraints. Finding an appropriate receiver-donor
102 relation is important. It can for instance be advantageous to choose it in such a way that
103 all receiver-donor pairs contain enough mutual information, which facilitates a meaningful
104 comparison. On the other hand, if a receiver-donor relation is chosen too restrictive, the set
105 of possible donors might be too small and the imputation problem is rendered infeasible. A
106 compromise solution in this case considers receiver-donor pairs where the missing values of
107 both the donor and receiver are at least mutually complementary.

108 If a data set has few missing values and more respondents than non-respondents, one
109 can often define a receiver-donor relation by considering only respondents as donors and non-
110 respondents as receivers. This relation ensures that the quality of each imputation can be
111 evaluated consistently. We call this receiver-donor relation the *respondent-non-respondent*
112 *relation*. Formally, the relation is defined by setting $u \prec v$ if and only if u is a non-respondent
113 and v is a respondent.

114 Possible receiver-donor pairs become comparable by introducing a distance.

115 **Definition 2.3.** *Let \mathcal{S} be a sample and \prec a receiver-donor relation. We call a symmetric*
116 *mapping*

$$117 \quad d : A_{\prec}(\mathcal{S}) \cup \{(u, v) : (v, u) \in A_{\prec}(\mathcal{S})\} \rightarrow \mathbb{R}_+$$

118 *a distance of \mathcal{S} .*

119 Common distances in the missing data context are the Gower distance [10], the Maha-
120 lanobis distance [23], and distances based on the *predictive mean* as discussed by [2]. As we
121 use the Gower distance for the computations in Section 5, we review its definition in A.

122 **3. Formulation of the Optimization Problem.** We have provided the necessary frame-
123 work for the formulation of the imputation problem for a given sample \mathcal{S} and a receiver-donor
124 relation \prec as a mathematical optimization problem. The objective is the minimization of the
125 sum over all distances of all imputation pairs, while respecting an upper bound on each donor's
126 maximal contribution. The upper bound for each donor v is denoted by $b_v \in \mathbb{Z}_+$ and the vec-
127 tor of all upper bounds by b . Usually, all donations are bounded by the same value, but our
128 model allows to choose individual donation constraints.

129 We model the problem as a graph problem by considering the union of the sets R and D ,
130 defined in (2.1) and (2.2), as the vertex set V of a graph $G = (V, E)$. Although in general R
131 and D are not disjoint, we simplify the construction without loss of generality in such a way
132 that each element in R and D is treated as a unique entity and thus as a single node in the
133 graph. Formally, this is achieved by extending the units with an $(q + 1)$ -th binary variable
134 indicating if the unit is in the set R or D and taking the union of the extended units. That
135 means a receiver unit u is replaced by an auxiliary unit $(u, 1)$ and a donor unit v is replaced
136 by $(v, 0)$. If u is a receiver and a donor, then we use two auxiliary units $(u, 1)$ and $(u, 0)$ to
137 construct the sets R and D .

138 We write $V = R \cup D$ for the disjoint union and obtain $|V| = |R| + |D|$, where $|\cdot|$ is the
139 cardinality of a set. The set of edges E of the graph is defined by $E := \{\{u, v\} \in \mathcal{S}^2 : u \prec v\}$.
140 Following the common notation, we let $m := |E|$ and $n := |V|$.

141 The graph defined above is by construction bipartite, i.e., there are only edges connecting
142 an element of R with an element of D but not between two elements of the same set. We
143 complete the construction by assigning distances d to each edge and an upper bound $b_v \in \mathbb{Z}_+$
144 to each donor $v \in D$.

145 Now, the imputation problem corresponds to the problem of finding a subset M of E in
146 which each $u \in R$ is covered by an edge of M , each $v \in D$ is an end-node of maximally $b_v \in \mathbb{Z}_+$
147 edges and M is minimal with respect to the distance d . Choosing such an M is achieved by
148 introducing indicator variables x , i.e., $x_e = 1$ if and only if the edge e is in M . This yields the
149 linear integer optimization problem (HDI1).

$$\begin{aligned}
& \min && \sum_{e \in E} d(e) x_e \\
& \text{s. t.} && \sum_{e \in \delta(u)} x_e = 1 && \forall u \in R \\
150 \text{ (HDI1)} &&& \sum_{e \in \delta(v)} x_e \leq b_v && \forall v \in D \\
&&& x_e \in \{0, 1\} && \forall e \in E,
\end{aligned}$$

151 where $\delta(u)$ denotes the set of edges covering u . Throughout our presentation, we assume that
152 the instances of Problem (HDI1) and its equivalent reformulations have at least one feasible
153 solution.

154 We first consider the case $b = \mathbf{1}$. The resulting problem loosely resembles a perfect mini-
155 mum weight matching. However, since in practice we have more possible donors than receivers,

156 a perfect matching does not exist. Therefore, rather than requiring a perfect matching, we
 157 compute a minimum-weight matching M , such that M is a maximum-cardinality matching.

158 In [29, Chapter 17] a modified Hungarian method for solving the perfect minimum-weight
 159 matching problem is proposed. Examining the method closely shows that the algorithm stops
 160 when a maximum-cardinality matching is found. Furthermore, such a matching is extreme
 161 [29, Theorem 17.2], i.e., it has the minimum weight under all matchings with the same size.
 162 Thus, the Hungarian method can be used to solve (HDI1) when fixing $b = \mathbf{1}$. As a matter
 163 of fact, Schrijver’s implementation of the Hungarian method can be regarded as a successive
 164 shortest path algorithm on a minimum-cost flow problem.

165 Although this observation only holds for the standard matching case $b = \mathbf{1}$, it is our goal
 166 to extend the described idea to arbitrary b . As seen in Section 4, an analogous approach is
 167 indeed applicable for the general case and the successive shortest path algorithm turns out to
 168 be a suitable choice for solving (HDI1), see Section 5.

169 We formulate the problem as a maximum-weighted b -matching problem, as defined by [29,
 170 Chapter 31]. This allows a broader perspective for theoretical examination, available solution
 171 methods, and problem complexities.

172 The transformation is straightforward and follows the idea in [25, pp. 610], see also B.1.
 173 We obtain the maximum-weighted b -matching problem (HDI2) as an equivalent formulation
 174 for (HDI1).

$$\begin{aligned}
 & \max \quad \sum_{e \in E} w(e) x_e \\
 & \text{s. t.} \quad \sum_{e \in \delta(u)} x_e \leq 1 \quad \forall u \in R \\
 & \quad \quad \sum_{e \in \delta(v)} x_e \leq b_v \quad \forall v \in D \\
 & \quad \quad x_e \in \mathbb{Z}_+ \quad \forall e \in E,
 \end{aligned}$$

175 (HDI2)

176 with $w(e) := |R| \max_{f \in E} d(f) + 1 - d(e)$.

177 **4. Solving the Nearest Neighbor Hot Deck Imputation Problem.** A straightforward
 178 approach for the solution of (HDI2), and equivalently (HDI1), is to find a solution for the LP
 179 relaxation of the problem. Indeed, it is well-known that the node-edge incidence matrix of a
 180 bipartite graph, like in (HDI2), is totally unimodular. Therefore, every optimal solution of
 181 the linear relaxation of (HDI2) is an integer point. The proof for the above classical result is
 182 for instance given in Chapter 3 of [4].

183 When dealing with small values for b_v , one could also introduce b_v copies of each donor node
 184 $v \in D$ and solve a weighted matching problem on the resulting graph. However, this is not a

185 polynomial time transformation and solving such a matching problem becomes computational
 186 expensive, even for small donor limits.

187 We are interested in developing more adapted solvers for the imputation problem. As it
 188 is equivalent to a maximum-weighted b -matching problem, there exist efficient algorithms to
 189 solve it. The algorithm presented in [8] takes $\mathcal{O}(m^2 \log n \log \|b\|_\infty)$ steps to solve a weighted b -
 190 matching problem. However, as the expression m^2 dominates this bound in dense graphs that
 191 typically occur in nearest neighbor hot deck imputation problems, it is preferable to reduce the
 192 worst-case number of iterations over all edges. This is possible by transforming (HD11) to a
 193 capacitated minimum cost flow problem. In practical applications, the cost scaling algorithm,
 194 developed independently by [28] and [5], and the network simplex algorithm are the most
 195 efficient algorithms to solve minimum cost flow problems [20]. The cost scaling algorithm was
 196 further enhanced by [9] to a runtime of $\mathcal{O}(nm \log(n^2/m) \log(nC))$ and the network simplex
 197 algorithm can find an optimal solution in $\mathcal{O}(nm \log(n) \log(nC))$ with $C = \|d\|_\infty$, see [31]. The
 198 disadvantage of both approaches is that they require integer distances, which is generally not
 199 the case for imputation problems. Thus, before applying any of the two algorithms, it is
 200 required to multiply the distances by a large factor and round them to the next integer. This
 201 has negative effects on the runtime because the logarithm of the maximal distance is a factor
 202 in the runtime bound.

203 We propose using the successive shortest path algorithm on the imputation problem pub-
 204 lished independently by [15], [14] and [6]. Although it generally runs only in pseudo-polynomial
 205 time, we show that it performs better on the imputation problem and reaches a worst-case
 206 runtime of $\mathcal{O}(|R|(m + n \log n))$. Moreover, the algorithm permits non-integral distances. As
 207 noted before, the successive shortest path approach is also in line with methods for bipartite
 208 perfect minimum-cost matchings, which further motivates its use.

209 **4.1. Transformation to a capacitated minimum cost flow problem.** Following a stan-
 210 dard approach for solving bipartite b -matching problems, we create a new source node s that
 211 contains all donations and transform the imputation problem into a capacitated minimum
 212 cost flow problem. This transformation is possible because the imputation graph is bipartite
 213 and fails for general graphs.

214 The directed graph (digraph) with vertex set $V_1 := R \cup D \cup \{s\}$ and arc set $A_1 :=$
 215 $A_\prec(\mathcal{S}) \cup \{(s, v) : v \in D\}$ is constructed. Furthermore, we define for $v \in D$ and $u \in R$ the
 216 weights

$$217 \quad d_1((v, u)) := \begin{cases} d(\{v, u\}), & \text{if } (v, u) \in A_\prec(\mathcal{S}), \\ 0, & \text{otherwise.} \end{cases}$$

218 On this digraph we can set up the minimum cost flow problem.

$$\begin{aligned}
& \min && \sum_{a \in A_1} d_1(a)x_a \\
& \text{s. t.} && \sum_{a \in \delta^{\text{out}}(v)} x_a - \sum_{a \in \delta^{\text{in}}(v)} x_a = 0 && \forall v \in D \\
219 \quad (\text{HDI3}) &&& - \sum_{a \in \delta^{\text{in}}(u)} x_a = -1 && \forall u \in R \\
&&& \sum_{a \in \delta^{\text{out}}(s)} x_a = |R| \\
&&& x_{(v,u)} \leq 1 && \forall v \in D, u \in R \\
&&& x_{(s,v)} \leq b_v && \forall v \in D,
\end{aligned}$$

220 where $\delta^{\text{in}}(v)$ denotes the set of arcs ending in v and $\delta^{\text{out}}(v)$ denotes the set of arcs starting
221 in v .

222 In each iteration the successive shortest path algorithm reduces the node imbalance of
223 two vertices by at least 1 for each of the two vertices, see [1, Chapter 9] and definitions
224 therein. Since the total, absolute imbalance of (HDI3) is $2|R|$, it requires $|R|$ iterations. In
225 each iteration a shortest path is searched in time $\mathcal{O}(m + n \log n)$ using Dijkstra's algorithm
226 with Fibonacci heaps [29]. Therefore, the successive shortest path algorithm takes $\mathcal{O}(|R|(m +$
227 $n \log n))$ time to find an optimal solution of (HDI3).

228 **5. Computational Study.** We have modeled the generalized hot deck imputation problem
229 with donor constraints in a mathematically concise way and presented how it can be solved as
230 a combinatorial optimization problem. In this section, we show that the theory also translates
231 to powerful algorithms for practical problem solving.

232 **Algorithms.** Our implementations comprise the network simplex (NSX), cost-scaling (CS),
233 and successive shortest path (SSP) algorithm. For all implementations the open source C++
234 LEMON graph library [7] is used. While NSX and CS algorithms are provided by LEMON,
235 we implemented the SSP algorithm using LEMON's data structures and its implementation
236 of Dijkstra's algorithm. Since both the network simplex and the cost-scaling approach require
237 integer distances on the arcs, we multiply all distances with a sufficiently large factor, in our
238 case 10^8 , and round with the *round half up* rule.

239 We assess the quality of the above algorithms in comparison to the heuristic methods
240 matrix minimization (MMIN) and Vogel approximation method (VAM) available in the R
241 package *HotDeckImputation* [18]. The MMIN method is a greedy algorithm and Vogel's
242 approximation is a heuristic for transportation problems proposed by [27]. Both methods

243 are mentioned as adequate tools by [16]. The package *HotDeckImputation* also contains a
 244 network simplex implementation which, however, does not consistently find optimal solutions
 245 and every so often failed to terminate. To our knowledge, *HotDeckImputation* is the only
 246 available solver for the donor-limited hot deck imputation problem so far.

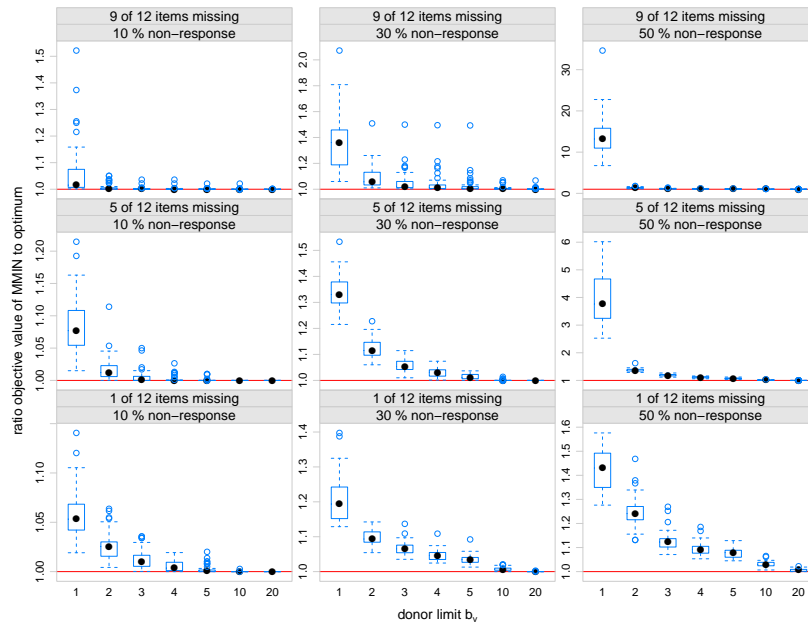


Figure 5.1. Ratios of the objective value of the matrix minimum method (MMIN) to the optimal value. Each boxplot consist of 49 independent samples of size 2000.

247 **Simulation setup.** The imputations are calculated on the synthetic data set AMELIA
 248 [24] using 12 variables of categorical, ordered, and continuous scale. We consider sample sizes
 249 of $\mathbf{size} = \{1000, 2000, 3000, 5000\}$ units. For each sample size, 49 simple random samples
 250 without replacement are drawn. Within of each of these samples, $\mathbf{unit} = \{10\%, 30\%, 50\%\}$
 251 of the units are randomly selected to contain missings. The number of missing items is then
 252 chosen to be $\mathbf{item} = \{1, 5, 9\}$ out of the total 12 items. The items to be missing are again
 253 drawn randomly. This can be regarded as an MCAR pattern, as discussed by [22]. In total,
 254 we create $|\mathbf{unit} \times \mathbf{item}| = 15$ scenarios for each of the $|\mathbf{size}| * 49 = 196$ samples. For each
 255 scenario in each sample the optimal solution for the donor limited hot deck nearest imputation
 256 problem is computed using the donor-limits $b_v = 1, \dots, 5, 10, 20$. The optimization problems
 257 are generated using the respondent-non-respondent relation and the Gower distance.

258 **Quality of solutions.** The imputation quality relies on the objective value and minimizing
 259 the sum over all distances implies improved statistical qualities of the imputed data set, see
 260 [16].

261 We compare the final objective values of MMIN with the optimal solution to the imputa-
 262 tion problem. In many samples the objective value of MMIN is much higher than the optimal
 263 value (Figure 5.1), up to a factor of 35. The approximation quality of MMIN is especially
 264 poor for the restrictive situation with a low donor-limit of $b_v = 1$. The greedy approximation
 265 becomes better for higher donor limits. This is plausible since donor limits pose no restriction
 266 if chosen sufficiently large. Thus, a greedy method finds an optimal solution in this case.
 267 The same patterns can be observed in Table 5.1 when comparing median objective values.
 268 Approximation quality improves when increasing b_v , but is very poor for $b_v = 1$. For limits
 269 above 5 the approximation quality of MMIN becomes tolerable.

b_v	1	2	3	4	5	10	20
CS	4.70	2.99	2.75	2.66	2.63	2.62	2.62
SSP	4.70	2.99	2.75	2.66	2.63	2.62	2.62
NSX	4.70	2.99	2.75	2.66	2.63	2.62	2.62
MMIN	19.88	4.05	3.29	2.94	2.79	2.64	2.63
VAM	16.01	-	-	-	-	-	-

Table 5.1

Median objective value for the scenario with 50% missing units, 5 items missing out 12, and sample size 2000.

270 Likewise, we compare the objective values produced by VAM to the optimal values in
 271 Table 5.1. The median of the objective values returned by VAM for the case $b_v = 1$ is three
 272 times higher than the median of the optimal values. We can not compare VAM to larger
 273 donor limits since the implementation provided in *HotDeckImputation* fails on all instances
 274 with $b_v > 1$ and non-respondents rates of less than 50%, i.e., if a perfect b -matching does
 275 not exist. We suspect that the code has to be adapted to handle excessive supply in the
 276 transformation of the transportation problem. In comparison to the MMIN method, VAM
 277 has a better approximation quality, which was also observed by [16].

278 In conclusion, both methods provided by *HotDeckImputation* do generally not find the
 279 optimum and return objective values which can be significantly worse than optimal value. The
 280 key assumption of the nearest neighbor hot deck imputation is that smaller distances imply
 281 an improved imputation quality. Thus, lower objective values lead to enhanced estimations
 282 of missing values and optimal solvers should be used whenever available. Furthermore, VAM
 283 is a rather uncommon method and to our understanding does not yield any advantages over

284 more sophisticated algorithms. Moreover, approximation qualities for this method are not
 285 guaranteed, which makes its usage even more questionable.

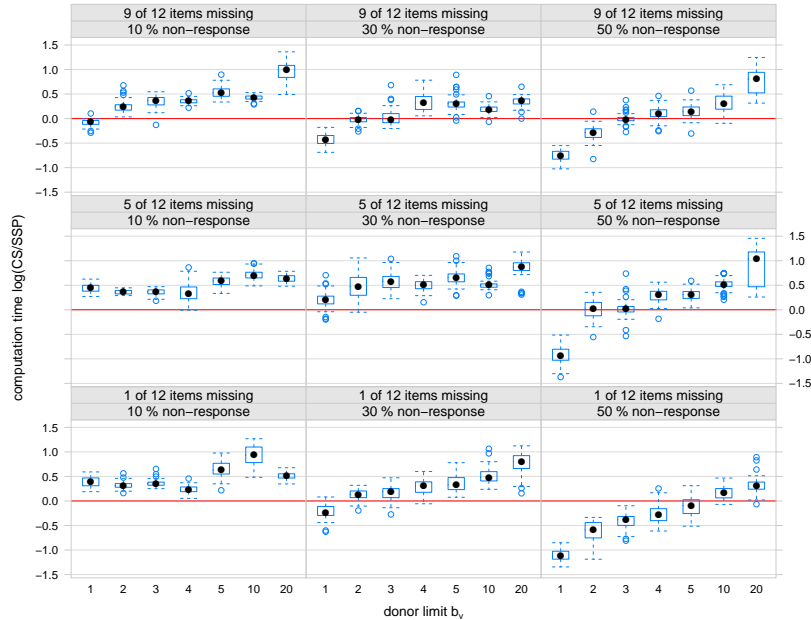


Figure 5.2. Log Ratios of measured computation times between the implementation of a successive shortest path algorithm (SSP) and the LEMON cost-scaling algorithm (CS). Each boxplot consist of 49 independent samples of size 2000.

286 **Computation times.** In Figure 5.2 the log ratio of the computation times of the CS to
 287 the SSP algorithm is depicted for the sample size 2000. The computation times of the CS
 288 algorithm are higher (non-negative log ratio) in most cases. However, if the imputation problem
 289 is restrictive, i.e., a high percentage of units are non-respondents and the donor limit is low,
 290 CS performs faster. This behavior is more pronounced for larger sample sizes (see Figure
 291 5.3), where SSP is never faster in the scenario with 50% missings and 1 of 12 items missing.
 292 Moreover, it can be observed that the performance of NSX is very similar to the performance
 293 of CS. Nevertheless, SSP is the superior choice in most settings, whereas CS or NSX should
 294 be considered when the imputation problem has few donors or is overly restrictive. For the
 295 scenario of sample size 2000, 50% missing units, and 5 out of 12 items missing we present the
 296 median computation times for the five algorithms SSP, CS, NSX, VAM, and MMIN in Table
 297 5.2.

298 Use of inexact heuristics might be justified when exact algorithms turn out to be too
 299 resource-consuming. However, we could not observe a considerable performance advantage
 300 for the *HotDeckImputation* implementations in Figure 5.4 and Table 5.2. In conclusion, the

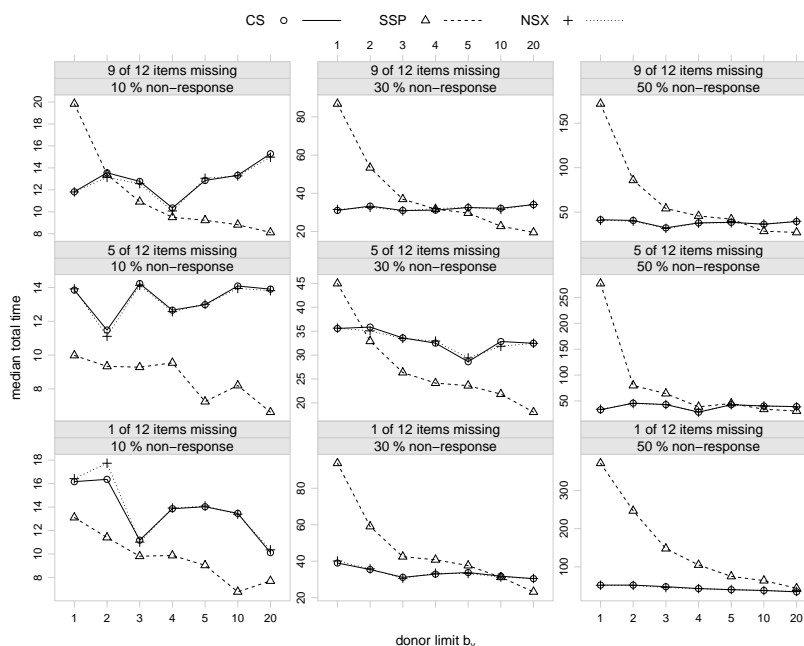


Figure 5.3. Median computation time in seconds of the LEMON cost-scaling algorithm (CS), the successive shortest path algorithm (SSP), and the network simplex (NSX). Each point is the median computation time over the 49 independent samples of size 5000.

301 heuristic methods do neither provide a sufficient approximation quality nor do they perform
 302 better in terms of running times.

b_v	1	2	3	4	5	10	20
CS	7.82	6.89	4.74	6.22	6.25	6.02	8.05
SSP	20.49	6.68	4.76	4.76	4.64	3.57	3.16
NSX	7.85	6.90	4.79	6.34	6.19	5.98	8.02
MMIN	8.96	9.80	10.73	10.47	11.05	8.97	13.45
VAM	8.01	-	-	-	-	-	-

Table 5.2

Median computation time for the scenario with 50% missing units, 5 items missing out of 12, and sample size 2000.

303 **6. Summary.** Nearest neighbor hot deck imputation is an important and commonly used
 304 tool in Data Science. [19] and [18] show that limiting the number of times a donor may be
 305 re-used can improve the imputation quality. However, the methodology for solving a donor-
 306 limited nearest neighbor hot deck imputation problem has not been examined thoroughly until
 307 now. We have closed this gap by modeling the imputation problem as a maximum-weighted

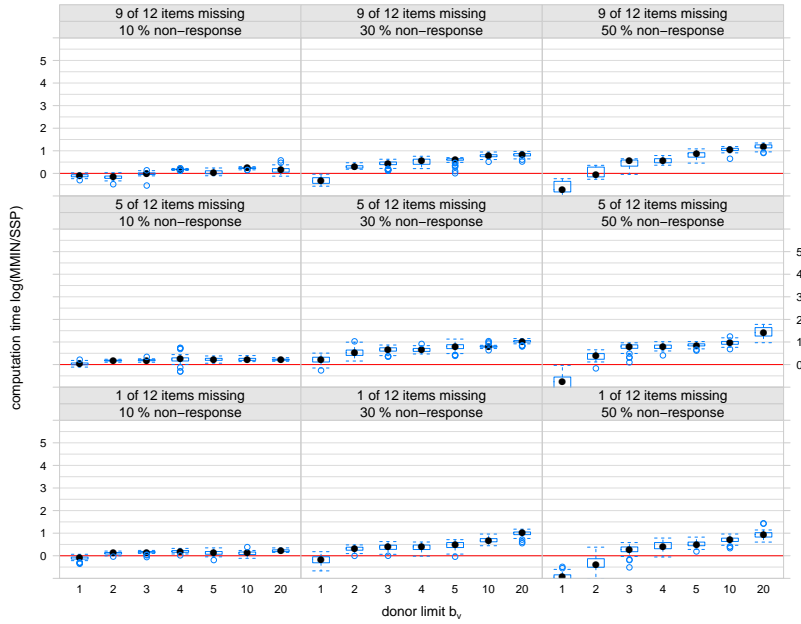


Figure 5.4. Log Ratios of measured computation times between the implementation of a successive shortest path algorithm (SSP) and the matrix minimum method (MMIN) from the HotDeckImputation R-Package. Each boxplot consist of 49 independent samples of size 2000.

308 b -matching, formulating a graph problem, and computing donor-limited nearest neighbor hot
 309 deck imputations with the help of exact combinatorial algorithms.

310 We have tested several algorithms on the open data set AMELIA for multiple scenarios.
 311 The computational study shows that our algorithms compute solutions that are superior to
 312 those of the priorly available solvers. Comparing the numerical performance of the algorithms,
 313 our implementation of the successive shortest path algorithm outperforms the other algorithms
 314 in most interesting cases. If the percentage of missing units is particularly high and donor
 315 limits are low, the usage of CS or NSX is more advantageous.

316 The presented graph theoretic approach gives both mathematically grounded, theoretical
 317 insights as well as fast algorithms for the computation of a global solution to the donor-limited
 318 nearest neighbor hot deck imputation problem in practice.

319 **Acknowledgements.** This research was supported within the research training group 2126
 320 *Algorithmic Optimization* (ALOP) founded by the German Research Foundation DFG.

321

REFERENCES

- 322 [1] R. K. AHUJA, T. L. MAGNANTI, AND J. B. ORLIN, *Network Flows: Theory, Algorithms and Applications*,
323 Prentice-Hall, Upper Saddle River, 1993.
- 324 [2] R. R. ANDRIDGE AND R. J. A. LITTLE, *A review of hot deck imputation for survey non-response*, Inter-
325 national Statistical Review, 78 (2010), pp. 40–64.
- 326 [3] M. BARNIER, *Edit and imputation for the 2006 Canadian census*, in Conference of European statisticians,
327 2005.
- 328 [4] D. BERTSIMAS AND R. WEISMANTEL, *Optimization over Integers*, Dynamic Ideas, Belmont, 2005.
- 329 [5] R. G. BLAND AND D. L. JENSEN, *On the computational behavior of a polynomial-time network flow*
330 *algorithm*, Mathematical Programming, 54 (1992), pp. 1–39.
- 331 [6] R. G. BUSACKER AND P. J. GOWEN, *A procedure for determining a family of minimum-cost network*
332 *flow patterns*, tech. report, Operations Research Office, Johns Hopkins University, 1960.
- 333 [7] B. DEZSŐ, A. JÜTTNER, AND P. KOVÁCS, *LEMON – An open source C++ graph template library*,
334 Electronic Notes in Theoretical Computer Science, 264 (2011), pp. 23–45. See also [http://lemon.cs.
335 elte.hu/trac/lemon](http://lemon.cs.elte.hu/trac/lemon).
- 336 [8] H. N. GABOW, *An efficient reduction technique for degree-constrained subgraph and bidirected network*
337 *flow problems*, in Proceedings of the Fifteenth Annual ACM Symposium on Theory of Computing,
338 1983, pp. 448–456.
- 339 [9] A. GOLDBERG AND R. TARJAN, *Solving minimum-cost flow problems by successive approximation*, in
340 Proceedings of the Nineteenth Annual ACM Symposium on Theory of Computing, 1987, pp. 7–18.
- 341 [10] J. C. GOWER, *A general coefficient of similarity and some of its properties*, Biometrics, 27 (1971),
342 pp. 857–871.
- 343 [11] S. GRUNWALD AND A. KRAUSE, *Umgang mit fehlenden Angaben in der Gebäude- und Wohnungszählung*
344 *2011*, Statistisches Bundesamt, Wirtschaft und Statistik, (2014).
- 345 [12] D. HAZIZA AND G. KUROMI, *Handling item non response in surveys*, Case Studies in Business, Industry
346 And Government Statistics, 1 (2014), pp. 102–118.
- 347 [13] D. HAZIZA AND É. LESAGE, *A discussion of weighting procedures for unit nonresponse*, Journal of Official
348 Statistics, 32 (2016), pp. 129–145.
- 349 [14] M. IRI, *A new method of solving transportation-network problems*, Journal of the Operations Research
350 Society of Japan, 3 (1960), pp. 27–87.
- 351 [15] W. S. JEWELL, *New methods in mathematical programming – optimal flow through networks with gains*,
352 Operations Research, 10 (1962), pp. 476–499.
- 353 [16] D. W. JOENSSEN, *Donor limited hot deck imputation: A constrained optimization problem*, in Studies
354 in Classification, Data Analysis, and Knowledge Organization, vol. 48, Springer, Berlin, Heidelberg,
355 2015, pp. 319–328.
- 356 [17] D. W. JOENSSEN, *Hot-Deck-Verfahren zur Imputation fehlender Daten: Auswirkungen des Donor-Limits*,
357 PhD thesis, Technische Universität Ilmenau, 2015.
- 358 [18] D. W. JOENSSEN AND U. BANKHOFER, *Donor limited hot deck imputation: effects on parameter estima-*
359 *tion*, Journal of Theoretical and Applied Computer Science, 6 (2012), pp. 58–70.
- 360 [19] G. KALTON AND L. KISH, *Two efficient random imputation procedures*, in Proceedings of the Survey
361 Research Methods Section, 1981, pp. 146–151.
- 362 [20] P. KOVÁCS, *Minimum-cost flow algorithms: an experimental evaluation*, Optimization Methods and
363 Software, 30 (2015), pp. 94–127.
- 364 [21] A. KOWARIK AND M. TEMPL, *Imputation with the R package VIM*, Journal of Statistical Software, 74
365 (2016), pp. 1–16.

- 366 [22] R. J. A. LITTLE AND D. B. RUBIN, *Statistical Analysis with Missing Data*, John Wiley & Sons, Hoboken,
367 2nd ed., 2002.
- 368 [23] P. C. MAHALANOBIS, *On the generalized distance in statistics*, Proceedings of the National Institute of
369 Sciences (Calcutta), 2 (1936), pp. 49–55.
- 370 [24] H. MERKLE, J. P. BURGARD, AND R. MÜNNICH, *The AMELIA Dataset – A Synthetic Universe for*
371 *Reproducible Research*, in Deliverable 23.1: Case Studies, Y. G. Berger, J. P. Burgard, A. Byrne,
372 A. Cernat, C. Giusti, P. Koksel, S. Lenau, S. Marchetti, H. Merkle, R. Münnich, I. Permanyer,
373 M. Pratesi, N. Salvati, N. Shlomo, D. Smith, and N. Tzavidis, eds., vol. WP23 – D23.1, 2016.
- 374 [25] G. L. NEMHAUSER AND L. A. WOLSEY, *Integer and Combinatorial Optimization*, John Wiley & Sons,
375 New York, 1988.
- 376 [26] A. OLINSKY, S. CHEN, AND L. HARLOW, *The comparative efficiency of imputation methods for missing*
377 *data in structural equation modeling*, European Journal of Operational Research, 151 (2003), pp. 53–
378 79.
- 379 [27] N. V. REINFELD AND W. R. VOGEL, *Mathematical Programming*, Prentice-Hall, 1958.
- 380 [28] H. RÖCK, *Scaling techniques for minimal cost network flows*, in Discrete Structures and Algorithms –
381 Proceedings of the 5th Conference on Graphtheoretic Concepts in Computer Science, 1980, pp. 181–
382 191.
- 383 [29] A. SCHRIJVER, *Combinatorial Optimization: Polyhedra and Efficiency*, Springer, Berlin, 2003.
- 384 [30] J. STACK, *The analysis of survey data*, European Journal of Operational Research, 81 (1995), pp. 1–16.
- 385 [31] R. E. TARJAN, *Dynamic trees as search trees via euler tours, applied to the network simplex algorithm*,
386 *Mathematical Programming*, 78 (1997), pp. 169–177.
- 387 [32] J. ZHANG AND H. AYTUG, *Comparison of imputation methods for discriminant analysis with strategically*
388 *hidden data*, European Journal of Operational Research, 255 (2016), pp. 522–530.

389 **Appendix A. Gower Distance.** Following [10], a popular method to measure proximity
390 is the so-called *Gower distance*. It differentiates between nominal and ordinal continuous
391 attributes and measures normalized, absolute differences between single variables. The total
392 sum of these differences is then called the Gower distance. This idea is rendered precisely
393 in the next two definitions. While the first Definition introduces pointwise distances for
394 dichotomous (i), quantitative (ii), and qualitative (iii) variables, the subsequent definition
395 combines the pointwise values in the actual Gower distance. A *quantitative* variable can be
396 ordered while a *qualitative* variable is categorical. A *dichotomous* variable has two states,
397 the + state and the – state. From + unambiguous information can be deduced, whereas –
398 generally resembles several possible instances. For example the question “Are you married?”
399 does not specify the possibilities single, divorced or widowed when answered with “no”.

400 **Definition A.1.** Let \mathcal{S} be a sample with q columns and $u, v \in \mathcal{S}$. The pointwise distance
401 $g_j(u, v)$ between the variables of two units u, v is given by:

402) If the j -th variable is dichotomous, then

$$403 \quad g_j(u, v) := \begin{cases} 0, & \text{if } u_j = v_j = +, \\ 1, & \text{otherwise.} \end{cases}$$

404) If the j -th variable is quantitative, then

$$405 \quad g_j(u, v) := \begin{cases} 0, & \text{if } u_j = \text{NA} \text{ or } v_j = \text{NA}, \\ \frac{|u_j - v_j|}{(M_j - m_j)}, & \text{otherwise.} \end{cases}$$

406 where $M_j := \max\{w_j : w \in \mathcal{S}\}$ and $m_j := \min\{w_j : w \in \mathcal{S}\}$.

407) If the j -th variable is qualitative, then

$$408 \quad g_j(u, v) := \begin{cases} 0, & \text{if } u_j = v_j, \\ 1, & \text{otherwise.} \end{cases}$$

409 We say that two values u_j and v_j are *comparable* if none of them is missing, i.e., $u_j \neq \text{NA}$
410 and $v_j \neq \text{NA}$, or in the case of dichotomous values if at least one of the two is $+$. The pointwise
411 definitions above cover the cases when values are not comparable, but the associated distances
412 are not considered in the final summation.

413 **Definition A.2.** Let μ be defined by

$$414 \quad \mu_j(u, v) := \begin{cases} 1, & \text{if } u_j \text{ is comparable to } v_j, \\ 0, & \text{else.} \end{cases}$$

415 Then, the Gower distance g is defined by

$$416 \quad g(u, v) := \frac{\sum_{j=1}^q \mu_j(u, v) g_j(u, v)}{\sum_{j=1}^q \mu_j(u, v)}$$

417 for all u, v with $\mu_j(u, v) = 1$ for at least one j .

418 In short, the Gower distance is defined as the mean over all $g_j(u, v)$ for which the units
419 u, v are comparable. The results in Sections 4 and 5 are formulated in terms of the Gower
420 distance as it is widely used in practice. However, our results remain valid for all distance
421 choices.

422 **Appendix B. Transformation to b -matching.** Every perfect minimum weighted matching
 423 problem can be formulated as a maximum weight matching problem [25, pp. 610]. This idea
 424 can be extended straightforwardly.

425 **Proposition B.1.** *Let $G = (U_1 \cup U_2, E)$ be a bipartite graph with edge weights $w : E \rightarrow \mathbb{R}$,*
 426 *$A \in \mathbb{R}^{m \times n}$ If the problem*

$$\begin{aligned}
 & \min \quad \sum_{e \in E} -w(e)x_e \\
 & \text{s. t.} \quad \sum_{e \in \delta(u)} x_e = 1 \quad \forall u \in U_1 \\
 & \quad \quad Ax \leq b \\
 & \quad \quad x_e \in \{0, 1\} \quad \forall e \in E
 \end{aligned}$$

427 (B.1)

428 *is feasible, then every optimal solution of*

$$\begin{aligned}
 & \max \quad \sum_{e \in E} (w(e) + \theta)x_e \\
 & \text{s. t.} \quad \sum_{e \in \delta(u)} x_e \leq 1 \quad \forall u \in U_1 \\
 & \quad \quad Ax \leq b \\
 & \quad \quad x_e \in \{0, 1\} \quad \forall e \in E
 \end{aligned}$$

429 (B.2)

430 *is also an optimal solution of (B.1).*