

1        ***GSimp: A Gibbs sampler based left-censored missing value***  
2                    ***imputation approach for metabolomics studies***

3

4    Runmin Wei<sup>1,2,¶</sup>, Jingye Wang<sup>1,¶,\*</sup>, Erik Jia<sup>3</sup>, Tianlu Chen<sup>4</sup>, Yan Ni<sup>1</sup>, Wei Jia<sup>1</sup>

5

6    <sup>1</sup> University of Hawaii Cancer Center, Honolulu, Honolulu, USA

7    <sup>2</sup> Department of Molecular Biosciences and Bioengineering, University of Hawaii at

8    Manoa, Honolulu, Honolulu, USA

9    <sup>3</sup> Punahou School, Honolulu, Honolulu, USA

10   <sup>4</sup> Shanghai Key Laboratory of Diabetes Mellitus and Center for Translational

11   Medicine, Shanghai Jiao Tong University Affiliated Sixth People's Hospital, Shanghai,

12   China

13

14

15   \* Corresponding author

16   Email: [jingyew@hawaii.edu](mailto:jingyew@hawaii.edu) (JW)

17

18

19   <sup>¶</sup> These authors contributed equally to this work.

20

21

22

## 23 **Abstract**

24 Left-censored missing values commonly exist in targeted metabolomics datasets and  
25 can be considered as missing not at random (MNAR). Improper data processing  
26 procedures for missing values will cause adverse impacts on subsequent statistical  
27 analyses. However, few imputation methods have been developed and applied to the  
28 situation of MNAR in the field of metabolomics. Thus, a practical left-censored  
29 missing value imputation method is urgently needed. We have developed an iterative  
30 Gibbs sampler based left-censored missing value imputation approach (GSimp). We  
31 compared GSimp with other three imputation methods on two real-world targeted  
32 metabolomics datasets and one simulation dataset using our imputation evaluation  
33 pipeline. The results show that GSimp outperforms other imputation methods in terms  
34 of imputation accuracy, observation distribution, univariate and multivariate analyses,  
35 and statistical sensitivity. The R code for GSimp, evaluation pipeline, vignette,  
36 real-world and simulated targeted metabolomics datasets are available at:  
37 <https://github.com/WandeRum/GSimp>.

38

## 39 **Author summary**

40 Missing values caused by the limit of detection/quantification (LOD/LOQ) were  
41 widely observed in mass spectrometry (MS)-based targeted metabolomics studies and  
42 could be recognized as missing not at random (MNAR). MNAR leads to biased

43 parameter estimations and jeopardizes following statistical analyses in different  
44 aspects, such as distorting sample distribution, impairing statistical power, etc.  
45 Although a wide range of missing value imputation methods was developed for  
46 –omics studies, a limited number of methods was designed appropriately for the  
47 situation of MNAR currently. To alleviate problems caused by MNAR and facilitate  
48 targeted metabolomics studies, we developed a Gibbs sampler based missing value  
49 imputation approach, called GSimp, which is public-accessible on GitHub. And we  
50 compared our method with existing approaches using an imputation evaluation  
51 pipeline on real-world and simulated metabolomics datasets to demonstrate the  
52 superiority of our method from different perspectives.

53

## 54 **Introduction**

55 Missing values are commonly observed in mass spectrometry (MS) based  
56 metabolomics datasets. Many statistical methods require a complete dataset, which  
57 makes missing data an inevitable problem for subsequent data analysis. Generally,  
58 there are three types of missing values, missing not at random (MNAR), missing at  
59 random (MAR) and missing completely at random (MCAR) [1,2]. Unexpected  
60 missing values are considered as MCAR if they originate from random errors and  
61 stochastic fluctuations during the data acquisition process (e.g., incomplete  
62 derivatization or ionization). MAR assumes the probability of a variable being  
63 missing depends on other observed variables [1,2]. Thus, missing values due to

64 suboptimal data preprocessing, e.g., inaccurate peak detection and deconvolution of  
65 co-eluting compounds can be defined as MAR. Targeted metabolomics studies have  
66 been widely used for the accurate quantification of specific groups of metabolites.  
67 Due to the limit of compound quantifications (LOQ), missing values are usually  
68 caused by signal intensities lower than LOQ, also known as left-censored missing,  
69 which can be assigned to MNAR.

70 The processing of missing values has been developed and studied in MS data, which  
71 is an indispensable step in the metabolomics data processing pipeline [3]. One simple  
72 but naïve solution is the substitution of missing by determined values, such as zero,  
73 half of the minimum value (HM) or  $LOQ/c$  where  $c$  denotes a positive integer.  
74 Determined value substitutions, although commonly applied for dealing with missing  
75 values in metabolomics studies [4–6], can significantly affect the subsequent  
76 statistical analyses in different ways, e.g. underestimate variances of missing variables,  
77 decrease statistical power, fabricate pseudo-clusters among observations, etc. [1].  
78 Advanced statistical imputation methods have been developed for –omics studies, e.g.,  
79 k-nearest neighbors (kNN) imputation [7], singular value decomposition (SVD)  
80 imputation [8,9], random forest (RF) imputation [10]. Several metabolomics data  
81 analysis software tools provide different methods of dealing with missing values  
82 [11–15]. MetaboAnalyst [15–17], one widely used metabolomics analysis toolkit,  
83 provides Probabilistic PCA (PPCA), Bayesian PCA (BPCA) and SVD imputation.  
84 However, these methods are mainly aiming at imputing MCAR/MAR and not suitable

85 for the situation of MNAR. A limited number of approaches dealing with  
86 left-censored missing values were applied by researchers [18,19]. Quantile regression  
87 approach for left-censored missing (QRILC) imputes missing data using random  
88 draws from a truncated distribution with parameters estimated using quantile  
89 regression [20]. Although this imputation keeps the overall distribution of missing  
90 parts compared to determined value substitutions, it may produce random results  
91 since no more information is used for the prediction of missing parts. Another  
92 imputation method recently developed for MNAR is k-nearest neighbor truncation  
93 (kNN-TN) by Shah, et al. [21]. This approach applies Maximum Likelihood  
94 Estimators (MLE) for the means and standard deviations of missing variables based  
95 on truncated normal distribution. Then a Pearson correlation based kNN imputation  
96 method was implemented on standardized data. Although the author stated that  
97 kNN-TN could impute both MNAR and MAR, the imputed values were entirely  
98 dependent on the nearest neighbors while no constraint was placed upon the  
99 imputation. Thus, this approach might cause an overestimation of missing values.

100 To reduce adverse effects caused by missing values during metabolomics data  
101 analyses, we developed a left-censored missing value imputation framework, GSimp,  
102 where a prediction model was embedded in an iterative Gibbs sampler. We then  
103 compared GSimp with HM, QRILC, and kNN-TN on two real-world metabolomics  
104 datasets and one simulation dataset to demonstrate the advantages of GSimp  
105 regarding imputation accuracy, observation distribution, univariate analysis,

106 multivariate analysis and sensitivity. Our findings indicate that GSimp is a robust  
107 method to handle left-censored missing values in targeted metabolomics studies.

## 108 **Results**

### 109 **Gibbs sampler in GSimp**

110 A variable containing missing elements from FFA dataset was randomly selected to  
111 track the sequence of corresponding parameters and estimates across the first 500  
112 iterations out of a total of 2000 ( $100 \times 20$ ) iterations using GSimp. From Fig 1, we  
113 can observe that both fitted value  $\hat{y}$  and sample value  $\tilde{y}$  reach to the convergence after  
114 iterations and the standard deviation estimate  $\sigma$  drop to a steady state with small  
115 values. In addition, an upper constraint for the distribution of  $\tilde{y}$  indicated that it was  
116 drawn from a truncated normal distribution.

117

118 **Fig 1. Sequentially parameters updating in GSimp.** The first 500 iterations out of a  
119 total of 2000 ( $100 \times 20$ ) iterations using GSimp where  $\hat{y}$ ,  $\tilde{y}$  and  $\sigma$  represent fitted value,  
120 sample value and standard deviation correspondingly.

121

### 122 **Imputation comparisons**

123 We evaluated four different MNAR imputation/substitution methods on FFA, BA  
124 targeted metabolomics and simulation datasets. First, we measured the imputation  
125 performances using label-free approaches. SOR was used to measure the imputation

126 accuracy regarding the imputed values of each missing variable. From the upper panel  
127 of Fig 2, we can observe that GSimp has the best performance with the lowest SOR  
128 across all varying numbers of missing variables in both FFA and BA datasets. To  
129 measure the extent of imputation induced distortion on observation distributions, the  
130 PCA-Procrustes analysis was conducted between the original data and imputed data.  
131 The lower panel of Fig 2 shows that GSimp has the lowest Procrustes sum of squared  
132 errors compared to other methods, which means GSimp kept the overall observation  
133 distribution of original dataset with the least distortions.

134

135 **Fig 2. Evaluations of different imputation methods using unlabeled approaches.**

136 SOR on FFA dataset (upper left) and BA dataset (upper right) along with different  
137 numbers of missing variables based on four imputation methods: HM (red circle),  
138 QRILC (green triangle), GSimp (blue square), and kNN-TN (purple cross).  
139 PCA-Procrustes sum of squared errors on FFA dataset (lower left) and BA dataset  
140 (lower right) along with different numbers of missing variables based on four  
141 imputation methods: HM (red circle), QRILC (green triangle), GSimp (blue square),  
142 and kNN-TN (purple cross).

143

144 Then, we measured the imputation performances with binary labels provided. We  
145 compared the results of univariate and multivariate analyses for imputed and original

146 datasets. Since this is a case-control study, student's  $t$ -tests were applied for univariate  
147 analyses. Then we compared the results by calculating Pearson's correlation between  
148 log-transformed  $p$ -values calculated from imputed and original data for missing  
149 variables. Again, GSimp performs best with the highest correlations among four  
150 methods (upper panel of Fig 3) along with different numbers of missing variables, and  
151 it implies GSimp keeps the most biological variations regarding the univariate  
152 analyses results. For the multivariate analyses, we applied PLS-DA to distinguish the  
153 group differences. Similarly, we conducted PLS-Procrustes analysis while PLS was  
154 employed as a supervised dimension reduction technique. The lower panel of Fig 3  
155 demonstrates that GSimp preferably restores the original observation distribution with  
156 the lowest Procrustes sum of squared errors among four imputation methods.

157

158 **Fig 3. Evaluations of different imputation methods using labeled approaches.**

159 Pearson's correlation between log-transformed  $p$ -values of student's  $t$ -tests on FFA  
160 dataset (upper left) and BA dataset (upper right) along with different numbers of  
161 missing variables based on four imputation methods: HM (red circle), QRILC (green  
162 triangle), GSimp (blue square), and kNN-TN (purple cross). PLS-Procrustes sum of  
163 squared errors on FFA dataset (lower left) and BA dataset (lower right) along with  
164 different numbers of missing variables based on four imputation methods: HM (red  
165 circle), QRILC (green triangle), GSimp (blue square), and kNN-TN (purple cross).



166

167 On the simulation dataset, we compared QRILC, kNN-TN, and GSimp using same  
168 approaches. Consistent results were recognized (S1 Fig), and GSimp presents the best  
169 performances on the simulation dataset with the lowest SOR and  
170 PCA/PLS-Procrustes sum of squared errors and the highest correlation of univariate  
171 analysis results. Moreover, to examine the influences of statistical power using  
172 different imputation methods, we calculated *TPR* as the capacities to detect  
173 differential variables on different imputation datasets. Again, with both *p*-cutoff of  
174 0.05 and 0.01, GSimp shows the overall highest *TPR* over different missing numbers  
175 (Fig 4). This implies that GSimp impairs the sensitivity to the least extent among  
176 three methods, which is reasonable since GSimp also keeps the highest correlation of  
177 *p*-values in previous comparisons.

178

179 **Fig 4. Evaluations of different imputation methods using TPR for various**  
180 **p-cutoffs on simulation dataset.** TPR along with different numbers of missing  
181 variables based on three imputation methods: QRILC (green triangle), GSimp (blue  
182 square), and kNN-TN (purple cross) among different *p*-cutoff=0.05 (left panel), and  
183 0.01 (right panel).

184

## 185 **Discussion**

186 The purpose of this study is to develop a left-censored missing value imputation  
187 approach for targeted metabolomics data analysis. We evaluated GSimp with other  
188 three imputation methods (a.k.a kNN-TN, QRILC, and HM) and suggested that  
189 GSimp was superior to others using different evaluation methods. To illustrate the  
190 performance of GSimp, we randomly selected one variable containing missing values  
191 from FFA dataset (Fig 5) to compare the imputed values and original values.  
192 Although determined value substitution (e.g. HM) were widely used by researchers in  
193 the field of metabolomics, our results indicated that HM could severely distort the  
194 data distribution (upper left panel of Fig 5), thus impairing subsequent analyses. In  
195 comparison, QRILC kept the overall data distribution and variances (upper right panel  
196 of Fig 5). However, random values could be generated by this approach since QRILC  
197 imputes each missing variable independently without utilizing the predictive  
198 information from other variables. Statistical learning based method, kNN-TN, applied  
199 a correlation based kNN algorithm with parameters of missing variables estimated  
200 with truncated normal distributions. This method utilized the information of highly  
201 correlated variables of targeted missing variable, thus kept a linear trend between  
202 original values and imputed values. However, since no constraint was applied for the  
203 imputation, a right shift of missing part might occur, causing imputed values to  
204 exceed the truncation point (lower left panel of Fig 5). In contrast, GSimp utilized the  
205 predictive information of other variables by employing a prediction model and held a

206 truncated normal distribution for each missing element simultaneously, which ensured  
207 a favorable linear trend between imputed and original values as well as a reasonable  
208 bound for the imputed values (lower right panel of Fig 5).

209

210 **Fig 5. Comparisons of imputed values and original values on an example**  
211 **variable.** Scatter plots of imputed values (X-axis) and original values (Y-axis) on one  
212 example missing variable while non-missing elements represented as blue dots and  
213 missing elements as red dots based on four imputation methods: HM (upper left),  
214 QRILC (upper right), kNN-TN (lower left), and GSimp (lower right). Rug plots show  
215 the distributions of imputed values and original values.

216

217 In our approach, truncated normal distribution was used for the constraint of  
218 imputation results in Gibbs sampler steps. We applied the minimum observed value of  
219 missing variable as an informative upper truncation point and  $-\infty$  as a non-informative  
220 lower truncation point considering the situation of left-censored missing. Other values  
221 could also be applied in real-world metabolomics analyses, such as a known LOQ of a  
222 metabolite can be set as an upper truncation point. Additionally, when signal intensity  
223 of certain compound is larger than the upper limit of quantification range or saturation  
224 during instrument analysis, an informative lower truncation point could be  
225 correspondingly applied for the right-censored missing value. What's more, when

226 non-informative bounds for both upper and lower limits (e.g.,  $+\infty$ ,  $-\infty$ ) were applied,  
227 our GSimp could be extended to the situation of MCAR/MAR. With the flexible  
228 usage of upper and lower limits, our approach may provide a versatile and powerful  
229 imputation technique for different missing types. For other –omics datasets with  
230 missing values (especially MNAR), e.g. single cell RNA-sequencing data, we could  
231 also apply this method with few modifications of our default settings. Thus, it is  
232 worthy to evaluate our approach, GSimp, in other complex scenarios in the future.

233 Since GSimp employed an iterative Gibbs sampler method, a large number of  
234 iterations (*iters\_all*=20, *iters\_each*=100) are preferable for the convergence of  
235 parameters. However, as we tested on the simulation dataset with different number of  
236 iterations, a much less iterations (*iters\_all*=10, *iters\_each*=50) won't severely affect  
237 the imputation accuracy (S2 Fig). Among iterations for the whole data matrix, we  
238 applied a sequential imputation procedure for missing variables from the least number  
239 of missing values to the most. Such sequential approach improves imputation  
240 performances compared to parallel imputation approach.

241

## 242 **Materials and Methods**

### 243 **Diabetes datasets**

244 We employed datasets from a study of comparing serum metabolites between obese  
245 subjects with diabetes mellitus (N=70) and healthy controls (N=130) where N

246 represents the number of observations. Dataset 1: a total of 42 free fatty acids (FFAs)  
247 were identified and quantified in those participants in order to evaluate their FFA  
248 profiles [22]. Dataset 2: a total of 34 bile acids (BAs) were identified and quantified  
249 in a similar way using different analytical protocol [23].

250

### 251 **Simulation dataset**

252 For the simulation dataset, we first calculated the covariance matrix  $Cov$  based on the  
253 whole diabetes dataset ( $P=76$ ) where  $P$  represents the number of variables. Then we  
254 generated two separated data matrices with the same number of 80 observations from  
255 multivariate normal distributions, representing two different biological groups. For  
256 each data matrix, the sample mean of each variable was drawn from a normal  
257 distribution  $N(0, 0.5^2)$  and  $Cov$  was kept using SVD. Then, two data matrices were  
258 horizontally (column-wise) stacked together as a complete data matrix ( $N \times P=160 \times 76$ )  
259 so that group differences were simulated and covariance was kept.

260

### 261 **MNAR generation**

262 For two real-world targeted metabolomics datasets, we generated a series of MNAR  
263 datasets by using the missing proportion (number of missing variables/number of total  
264 variables) from 0.1 to 0.6 in a step of 0.05 with MNAR cut-off for each missing  
265 variable drawn from a uniform distribution  $U(0.1, 0.5)$  The elements lower than the  
266 corresponding cut-off were removed and replaced with NA. For the simulation dataset,

267 we generated a series of MNAR datasets by using the missing proportion from 0.1 to  
268 0.8 step by 0.1 with MNAR cut-off drawn from  $U(0.3, 0.6)$  for a more rigorous  
269 testing.

270

## 271 **Prediction model**

272 A prediction model was employed for the prediction of missing values by setting a  
273 targeted missing variable as outcome and other variables as predictors. Different  
274 prediction models, e.g., linear regression, elastic net [24], regression trees [25] and  
275 random forest [26], etc. could be embedded in our imputation framework. Elastic net  
276 was applied in our approach as an ideal prediction model considering its stability,  
277 accuracy, and efficiency. This model is a regularized regression with the combination  
278 of L1 and L2 penalties of the LASSO [27] and ridge [28] methods. The estimates of  
279 regression coefficients in elastic net are defined as

$$280 \quad \hat{\beta} = \operatorname{argmin}_{\beta} (\|y - X\beta\|^2 + \lambda[(1 - \alpha)/2\|\beta\|_2^2 + \alpha\|\beta\|_1]) \quad (1)$$

281 The L2 penalty  $(1 - \alpha)/2\|\beta\|_2^2$  improves the model's robustness by controlling the  
282 multicollinearities among variables which are widely existed in high-dimensional  
283 -omics data. And the L1 penalty  $\alpha\|\beta\|_1$  controls the number of predictors by  
284 assigning zero coefficients to the "unnecessary" predictors. From a Bayesian point of  
285 view, the regularization is a mixture of Gaussian and Laplacian prior distributions of  
286 coefficients which can pull the full model of maximum likelihood estimates  
287  $\|y - X\beta\|^2$  towards the null model of prior coefficients distribution, thus controls the

288 risk of overfitting and increase the model robustness. R package *glmnet* was used for  
289 the elastic net. We set hyperparameters  $\lambda$  as 0.01 (default setting for high-dimensional  
290 data) and  $\alpha$  as 0.5 (an equally mixture of LASSO and ridge penalties) [29].

291

## 292 **Gibbs sampler**

293 Gibbs sampler is a Markov Chain Monte Carlo (MCMC) technique that sequentially  
294 updates parameters while others are fixed. It can be used to generate posterior  
295 samples. For each missing variable in the dataset, we applied a Gibbs sampler to  
296 impute the missing values by sampling from a truncated normal distribution with  
297 prediction model fitted value as mean and root mean square deviation (RMSD) of  
298 missing part as standard deviation while truncated by specified cut-points. Assuming  
299 we have a  $n \times p$  data matrix  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \dots, \mathbf{X}_p)$  with only one variable  $\mathbf{X}_j$   
300 containing left-censored missing values. We denote  $\mathbf{X}_j$  as  $\mathbf{y}$  and the missing part as  $\mathbf{y}_m$   
301 with length  $m$  and non-missing part as  $\mathbf{y}_f$  with length  $f$ , and the rest of matrix  $\mathbf{X}_j$  as  $\mathbf{X}'$ .  
302 We can then set the lower truncation point  $lo$  as  $-\infty$  (centralized data) or 0 (original  
303 data) and upper  $hi$  as the minimum value of  $\mathbf{y}_f$  or a given LOQ. The truncation bounds  
304 ensure imputation results are constrained within  $[lo, hi]$ . Then, the Gibbs sampler  
305 approach can be described as following steps:

306 Step-1 (initialization): we initialize missing values (QRILC in our case), and get  $\mathbf{y}'$ ;

307 Step-2 (prediction): we then build a prediction model (elastic net in our case):  $\mathbf{y}' \sim \mathbf{X}'$ ;

308 Step-3 (estimation): based on the prediction model, we get the predicted value  $\hat{\mathbf{y}}$  and

309 the root mean square deviation (RMSD) of missing part  $\sigma = \sqrt{\frac{\sum_{i=1}^m (\hat{y}_{m_i} - y'_{m_i})^2}{m}}$  where

310  $y'_{m_i}$  and  $\hat{y}_{m_i}$  are  $i$ th initialized/imputed value and fitted value respectively;

311 Step-4 (sampling): we draw sample  $\tilde{y}_{m_i}$  from a truncated normal distribution

312  $N(\hat{y}_{m_i}, \sigma^2 \mid [lo, hi])$  for  $i$ th missing element and update  $\mathbf{y}'$ .

313 We iteratively repeat step-2 to step-4 and update  $\mathbf{X}_j$ .

314

### 315 **GSimp framework**

316 A whole data matrix  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \mathbf{X}_3, \dots, \mathbf{X}_p)$  contains a number of  $k$  ( $k \leq p$ )

317 left-censored missing variables. We present our imputation framework as following

318 algorithm.

---

**Algorithm:** Gibbs sampler based left-censored missing value imputation approach

---



---

**Require:**  $X$  an  $n \times p$  data matrix,  $iters\_all$  the number of iterations for imputing the whole matrix  $X$ ,  $iters\_each$  the number of iterations for imputing each missing variable, a vector of upper limits  $U$  ( $+\infty$  for non-missing variables) and a vector of lower limits  $L$  ( $-\infty$  for non-missing variables) with length  $p$ .

1.  $X^{imp} \leftarrow$  initialize the missing values for  $X$ ;
  2.  $K \leftarrow$  vector of indices of missing variables in  $X$  with increasing amount of missing values;
  3. **for** 1: $iters\_all$  **do**
  4.     **for**  $j$  in  $K$  **do**
  5.          $y' \leftarrow X_j^{imp}$ ,  $y'$  can be divided into two parts:  $y'_m$  is a vector of the imputed part (original missing part) with length  $m$  and  $y'_f$  is a vector of the non-missing part with length  $f$  while  $n = m + f$ ;
  6.          $X' \leftarrow X_{-j}^{imp}$ , represents the matrix  $X$  with  $j$ th column removed;
  7.          $lo \leftarrow L_j$  and  $hi \leftarrow U_j$ ;
  8.         **for** 1: $iters\_each$  **do**
  9.             Gibbs sampler step 2 to 4;
  10.         **end for**
  11.         Update  $X_j^{imp}$ ;
  12.     **end for**
  13. **end for**
  14. **return**  $X^{imp}$
- 

319

## 320 **Other imputation approaches**

321 Other three left-censored missing imputation/substitution methods were conducted in  
322 our study for performance comparison:

- 323 • kNN-TN (Truncation  $k$ -nearest neighbors imputation) [21]: this method applied a  
324 Newton-Raphson (NR) optimization to estimate the truncated mean and standard

325 deviation. Then, Pearson correlation was calculated based on standardized data  
326 followed by correlation-based kNN imputation.

327 • QRILC (Quantile Regression Imputation of Left-Censored data) [18,30]: this  
328 method imputes missing elements randomly drawing from a truncated distribution  
329 estimated by a quantile regression. R package *imputeLCMD* was applied for this  
330 imputation approach.

331 • HM (Half of the Minimum): This method replaces missing elements with half of  
332 the minimum of non-missing elements in the corresponding variable.

### 333 **Assessments of performance**

334 The assessments of imputation performance were conducted using an imputation  
335 evaluation pipeline from our previous study with both unlabeled and labeled  
336 measurements [31], which is accessible through:

337 <https://github.com/WandeRum/MVI-evaluation>. Unlabeled measurements include the  
338 NRMSE-based sum of ranks (SOR), principal component analysis (PCA)-Procrustes  
339 analysis while labeled measurements include correlation analysis for univariate results,  
340 partial least square (PLS)-Procrustes analysis. R package *vegan* was applied for  
341 Procrustes analysis [32] and *ropls* was applied for PLS analysis [33].

342 Furthermore, we evaluated the impacts of different imputation methods on the  
343 statistical sensitivity of detecting biological variances. On the simulation dataset, we  
344 calculated  $p$ -values from student's  $t$ -tests between two groups from original as well as  
345 imputed datasets. We marked a set  $S$  as real differential variables at a significant level

346 of  $p$ -cutoff (e.g. 0.05) from original simulation data, and a set  $S'$  as detected  
347 differential variables at the same significant level from imputed simulation data. Then  
348 we calculated the true positive rate  $TPR = \frac{\# of (S \cap S')}{\# of S}$  to evaluate the effects of  
349 different imputation methods in terms of detecting differential variables.

350

## 351 **Conclusion**

352 A practical left-censored missing value imputation method is needed in the field of  
353 metabolomics. We develop a new imputation approach GSimp that outperforms  
354 traditional determined value substitution method (HM) and other approaches (QRILC,  
355 and kNN-TN) for MNAR situations. GSimp utilized predictive information of  
356 variables and held a truncated normal distribution for each missing element  
357 simultaneously via embedding a prediction model into the Gibbs sampler framework.  
358 With proper modifications on the parameter settings, e.g. truncation points, GSimp  
359 may be applicable to handle different types of missing values and in different -omics  
360 studies, thus deserved to be further explored in the future.

361

362 **References**

- 363 1. Gelman A, Hill J. Data analysis using regression and multilevel/hierarchical  
364 models [Internet]. Cambridge University Press. 2006. doi:10.2277/0521867061
- 365 2. Little RJ a, Rubin DB. Statistical Analysis with Missing Data. Statistical  
366 analysis with missing data Second edition. 2002. doi:10.2307/1533221
- 367 3. Hrydziuszko O, Viant MR. Missing values in mass spectrometry based  
368 metabolomics: An undervalued step in the data processing pipeline.  
369 *Metabolomics*. 2012;8: 161–174. doi:10.1007/s11306-011-0366-4
- 370 4. Guo L, Milburn M V, Ryals JA, Lonergan SC, Mitchell MW, Wulff JE, et al.  
371 Plasma metabolomic profiles enhance precision medicine for volunteers of  
372 normal health. *Proc Natl Acad Sci*. 2015;112: E4901–E4910.  
373 doi:10.1073/pnas.1508425112
- 374 5. Liu J-J, Ghosh S, Kovalik J-P, Ching J, Choi HW, Tavintharan S, et al.  
375 Profiling of plasma metabolites suggests altered mitochondrial fuel usage and  
376 remodelling of sphingolipid metabolism in individuals with type 2 diabetes and  
377 kidney disease. *Kidney Int Reports*. 2016;2: 470–480.  
378 doi:10.1016/j.ekir.2016.12.003
- 379 6. Butte NF, Liu Y, Zakeri IF, Mohny RP, Mehta N, Voruganti VS, et al. Global  
380 metabolomic profiling targeting childhood obesity in the Hispanic population.  
381 *Am J Clin Nutr*. 2015;102: 256–267. doi:10.3945/ajcn.115.111872
- 382 7. Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R, et al.

- 383 Missing value estimation methods for DNA microarrays. *Bioinformatics*.  
384 2001;17: 520–525. doi:10.1093/bioinformatics/17.6.520
- 385 8. Hastie T, Tibshirani R, Sherlock G. Imputing missing data for gene expression  
386 arrays. Tech Report, Div Biostat Stanford Univ. 1999; 1–9.
- 387 9. Stacklies W, Redestig H, Scholz M, Walther D, Selbig J. pcaMethods - A  
388 bioconductor package providing PCA methods for incomplete data.  
389 *Bioinformatics*. 2007;23: 1164–1167. doi:10.1093/bioinformatics/btm069
- 390 10. Stekhoven DJ, Bühlmann P. Missforest-Non-parametric missing value  
391 imputation for mixed-type data. *Bioinformatics*. 2012;28: 112–118.  
392 doi:10.1093/bioinformatics/btr597
- 393 11. Mak TD, Laiakis EC, Goudarzi M, Fornace AJ. MetaboLyzer: A novel  
394 statistical workflow for analyzing postprocessed LC-MS metabolomics data.  
395 *Anal Chem*. 2014;86: 506–513. doi:10.1021/ac402477z
- 396 12. Katajamaa M, Miettinen J, Oresic M. MZmine: toolbox for processing and  
397 visualization of mass spectrometry based molecular profile data.  
398 *Bioinformatics*. 2006;22: 634–6. doi:10.1093/bioinformatics/btk039
- 399 13. Kessler N, Neuweger H, Bonte A, Langenkämper G, Niehaus K, Nattkemper  
400 TW, et al. MeltDB 2.0-advances of the metabolomics software system.  
401 *Bioinformatics*. 2013;29: 2452–2459. doi:10.1093/bioinformatics/btt414
- 402 14. Luedemann A, Von Malotky L, Erban A, Kopka J. TagFinder: Preprocessing  
403 software for the fingerprinting and the profiling of gas chromatography-mass

- 404 spectrometry based metabolome analyses. *Methods Mol Biol.* 2012;860:  
405 255–286. doi:10.1007/978-1-61779-594-7\_16
- 406 15. Xia J, Sinelnikov I V., Han B, Wishart DS. MetaboAnalyst 3.0-making  
407 metabolomics more meaningful. *Nucleic Acids Res.* 2015;43: W251–W257.  
408 doi:10.1093/nar/gkv380
- 409 16. Xia J, Psychogios N, Young N, Wishart DS. MetaboAnalyst: A web server for  
410 metabolomic data analysis and interpretation. *Nucleic Acids Res.* 2009;37.  
411 doi:10.1093/nar/gkp356
- 412 17. Xia J, Mandal R, Sinelnikov I V., Broadhurst D, Wishart DS. MetaboAnalyst  
413 2.0-a comprehensive server for metabolomic data analysis. *Nucleic Acids Res.*  
414 2012;40. doi:10.1093/nar/gks374
- 415 18. Lazar C, Gatto L, Ferro M, Bruley C, Burger T. Accounting for the Multiple  
416 Natures of Missing Values in Label-Free Quantitative Proteomics Data Sets to  
417 Compare Imputation Strategies. *J Proteome Res.* 2016;15: 1116–1125.  
418 doi:10.1021/acs.jproteome.5b00981
- 419 19. Shah JS, Rai SN, DeFilippis AP, Hill BG, Bhatnagar A, Brock GN. Distribution  
420 based nearest neighbor imputation for truncated high dimensional data with  
421 applications to pre-clinical and clinical metabolomics studies. *BMC*  
422 *Bioinformatics.* 2017;18: 114. doi:10.1186/s12859-017-1547-6
- 423 20. Lazar C, Gatto L, Ferro M, Bruley C, Burger T. Accounting for the Multiple  
424 Natures of Missing Values in Label-Free Quantitative Proteomics Data Sets to

- 425 Compare Imputation Strategies. *J Proteome Res.* 2016;15: 1116–1125.  
426 doi:10.1021/acs.jproteome.5b00981
- 427 21. Shah JS, Rai SN, DeFilippis AP, Hill BG, Bhatnagar A, Brock GN. Distribution  
428 based nearest neighbor imputation for truncated high dimensional data with  
429 applications to pre-clinical and clinical metabolomics studies. *BMC*  
430 *Bioinformatics.* 2017;18: 114. doi:10.1186/s12859-017-1547-6
- 431 22. Ni Y, Zhao L, Yu H, Ma X, Bao Y, Rajani C, et al. Circulating Unsaturated  
432 Fatty Acids Delineate the Metabolic Status of Obese Individuals.  
433 *EBioMedicine.* 2015;2: 1513–1522. doi:10.1016/j.ebiom.2015.09.004
- 434 23. Lei S, Huang F, Zhao A, Chen T, Chen W, Xie G, et al. The ratio of  
435 dihomog- $\gamma$ -linolenic acid to deoxycholic acid species is a potential biomarker  
436 for the metabolic abnormalities in obesity. *FASEB J. Federation of American*  
437 *Societies for Experimental Biology;* 2017; fj.201700055R.  
438 doi:10.1096/fj.201700055R
- 439 24. Zou H, Hastie T. Regularization and variable selection via the elastic net. *J R*  
440 *Stat Soc Ser B Stat Methodol.* 2005;67: 301–320.  
441 doi:10.1111/j.1467-9868.2005.00503.x
- 442 25. Breiman L, Friedman JH, Olshen RA, Stone CJ. Classification and Regression  
443 Trees. The Wadsworth statisticsprobability series. 1984.  
444 doi:10.1371/journal.pone.0015807
- 445 26. Breiman L. Random forests. *Mach Learn.* 2001;45: 5–32.

- 446 doi:10.1023/A:1010933404324
- 447 27. Tibshirani R. Regression Selection and Shrinkage via the Lasso [Internet].  
448 Journal of the Royal Statistical Society B. 1996. pp. 267–288.  
449 doi:10.2307/2346178
- 450 28. Hoerl AE, Kennard RW. Ridge Regression: Biased Estimation for  
451 Nonorthogonal Problems. Technometrics. 1970;12: 55–67.  
452 doi:10.1080/00401706.1970.10488634
- 453 29. Friedman AJ, Hastie T, Simon N, Tibshirani R, Hastie MT. Lasso and  
454 Elastic-Net Regularized Generalized Linear Models [Internet]. 2015. Available:  
455 <http://www.jstatsoft.org/v33/i01/>.
- 456 30. Lazar C. Imputation of left-censored missing data using QRILC method  
457 [Internet]. 2015.
- 458 31. Wei R, Wang J, Su M, Jia E, Chen T, Ni Y. Missing Value Imputation Approach  
459 for Mass Spectrometry-based Metabolomics Data. bioRxiv. 2017; Available:  
460 <http://biorxiv.org/content/early/2017/08/17/171967.abstract>
- 461 32. Oksanen J. Multivariate Analysis of Ecological Communities in R: vegan  
462 tutorial [Internet]. 2015.
- 463 33. Thévenot EA, Roux A, Xu Y, Ezan E, Junot C. Analysis of the Human Adult  
464 Urinary Metabolome Variations with Age, Body Mass Index, and Gender by  
465 Implementing a Comprehensive Workflow for Univariate and OPLS Statistical  
466 Analyses. J Proteome Res. 2015;14: 3322–3335.



467      doi:10.1021/acs.jproteome.5b00354

468

469 **Supporting information**

470

471 **S1 Fig. Evaluations of different imputation methods on simulation dataset.** SOR

472 (upper left), PCA-Procrustes sum of squared errors (upper right), Pearson's correlation

473 between log-transformed  $p$ -values of student's  $t$ -tests (lower left), and PLS-Procrustes

474 sum of squared errors (lower right) on simulation dataset along with different

475 numbers of missing variables based on three imputation methods: QRILC (green

476 triangle), GSimp (blue square), and kNN-TN (purple cross).

477

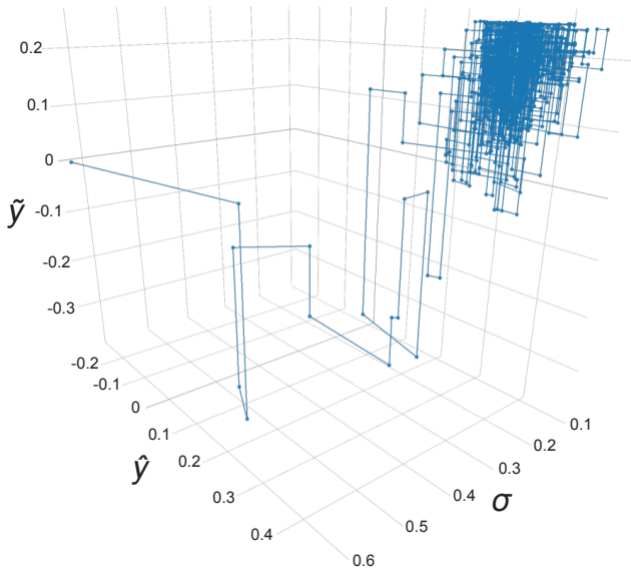
478 **S2 Fig. Evaluations of different numbers of iterations using GSimp on simulation**

479 **dataset.** SOR on simulation dataset along with different numbers of missing variables

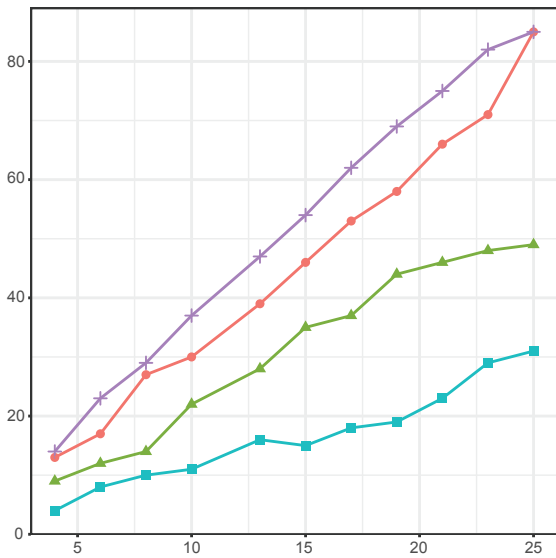
480 based on four different numbers of iterations:  $iters\_each=50$  and  $iters\_all=20$  (red

481 circle),  $iters\_each=100$  and  $iters\_all=20$  (green triangle),  $iters\_each=50$  and

482  $iters\_all=10$  (blue square),  $iters\_each=100$  and  $iters\_all=10$  (purple cross).

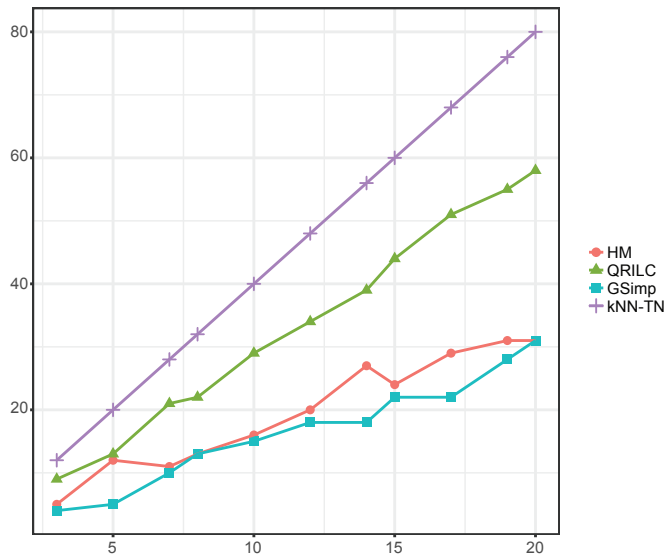


FFA



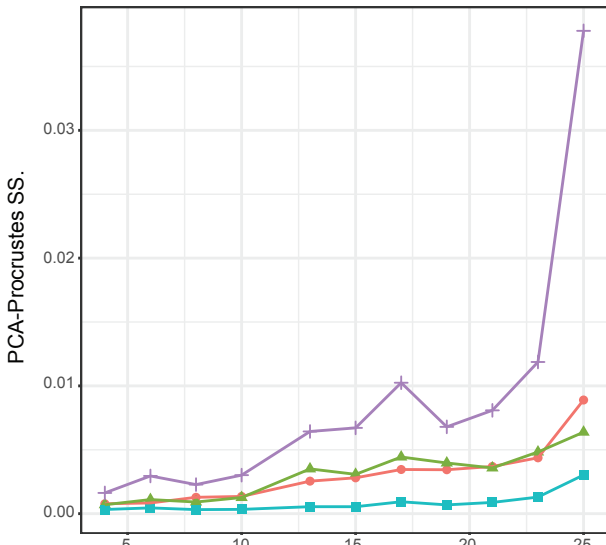
Number of missing variables

BA



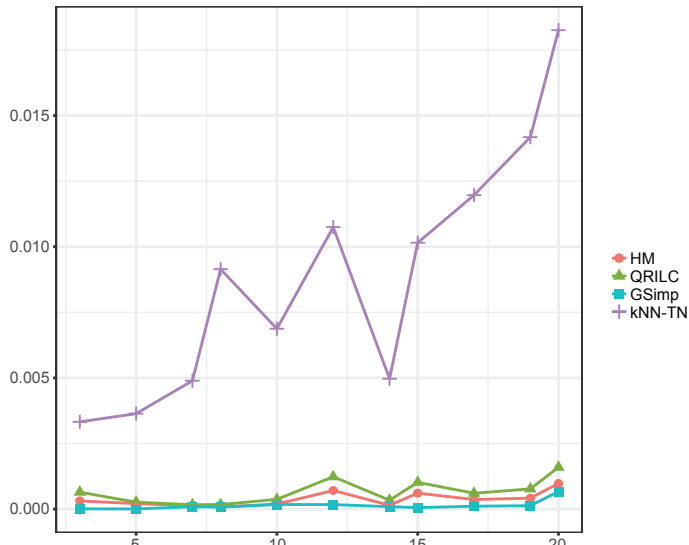
◆ HM  
▲ QRILC  
■ GSimp  
+ kNN-TN

FFA



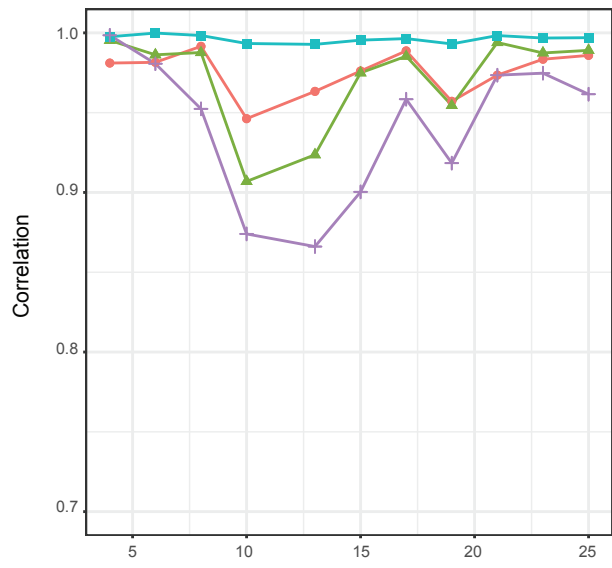
Number of missing variables

BA

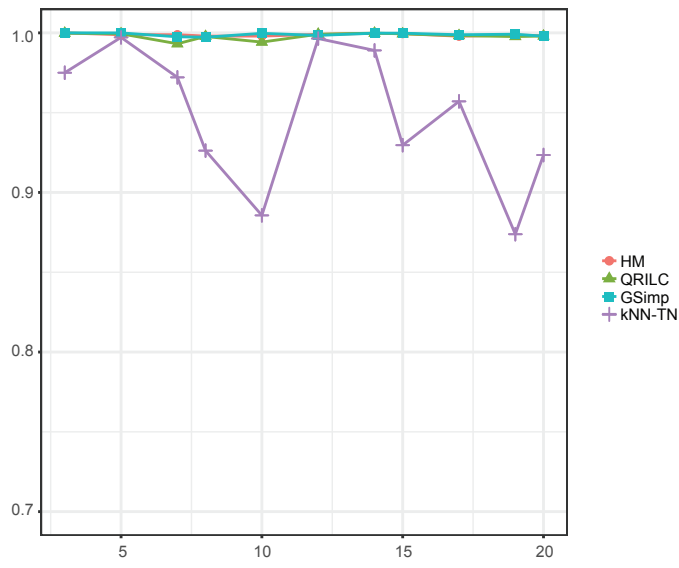


◆ HM  
▲ QRILC  
■ GSimp  
+ kNN-TN

FFA



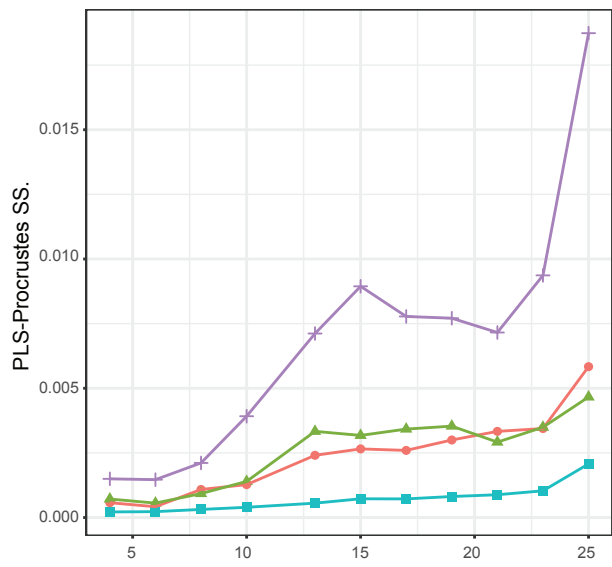
BA



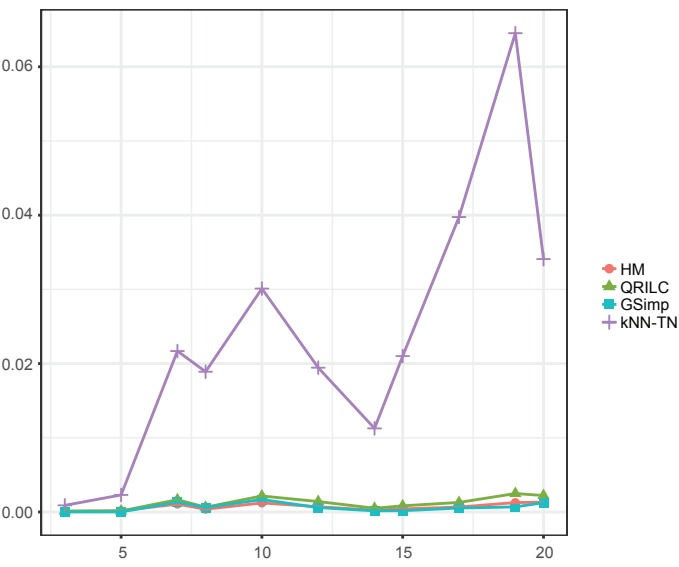
● HM  
▲ QRILC  
■ GSimp  
+ kNN-TN

Number of missing variables

FFA



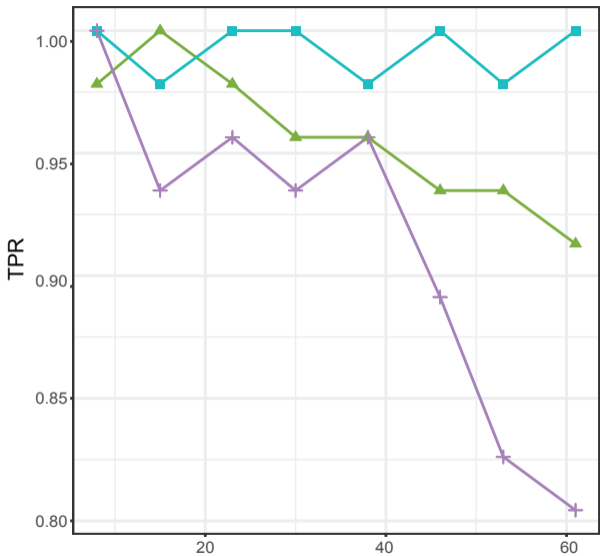
BA



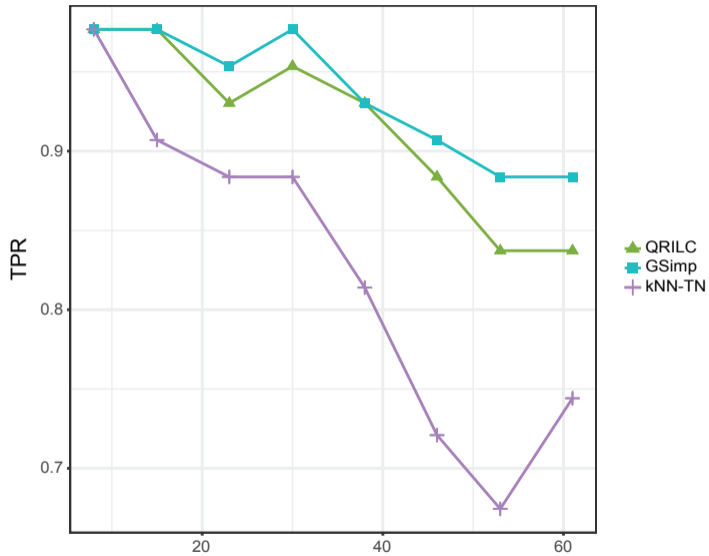
● HM  
▲ QRILC  
■ GSimp  
+ kNN-TN

Number of missing variables

p-cutoff=0.05

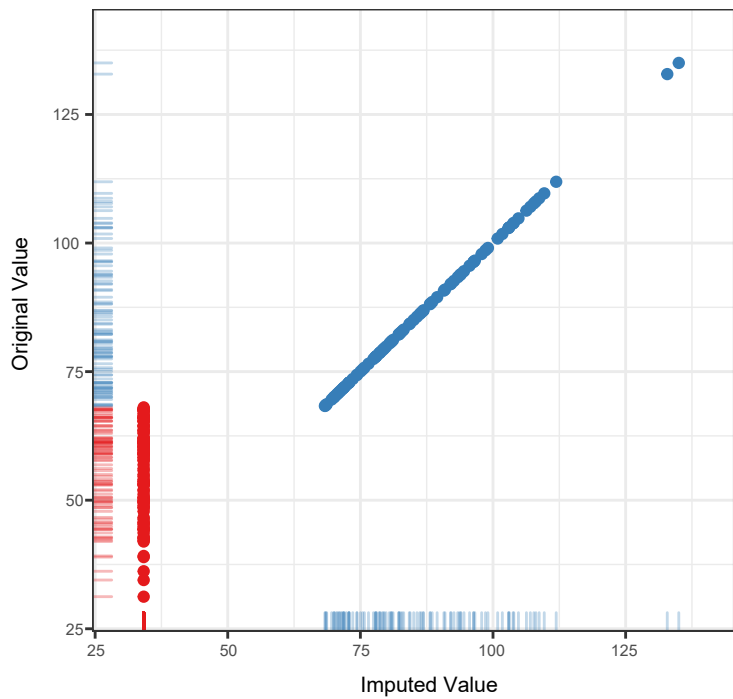


p-cutoff=0.01

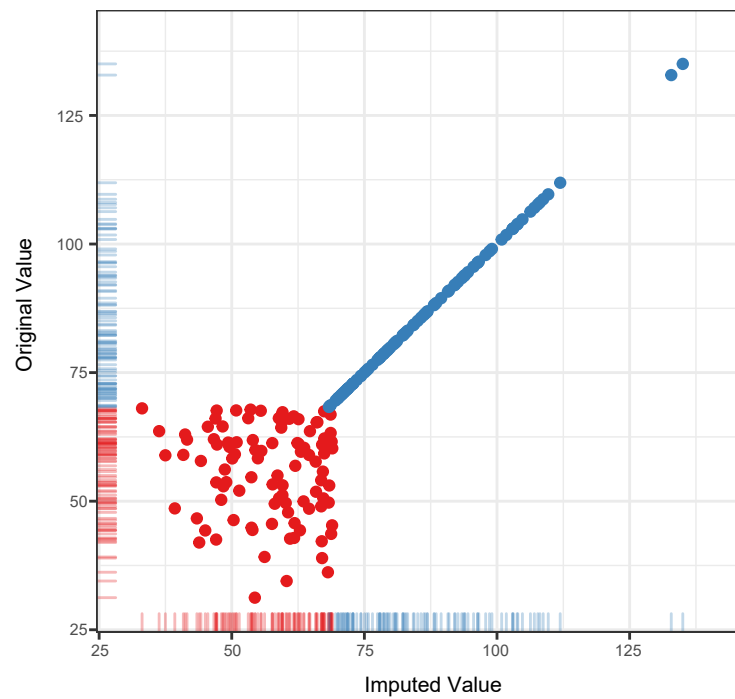


Number of missing variables

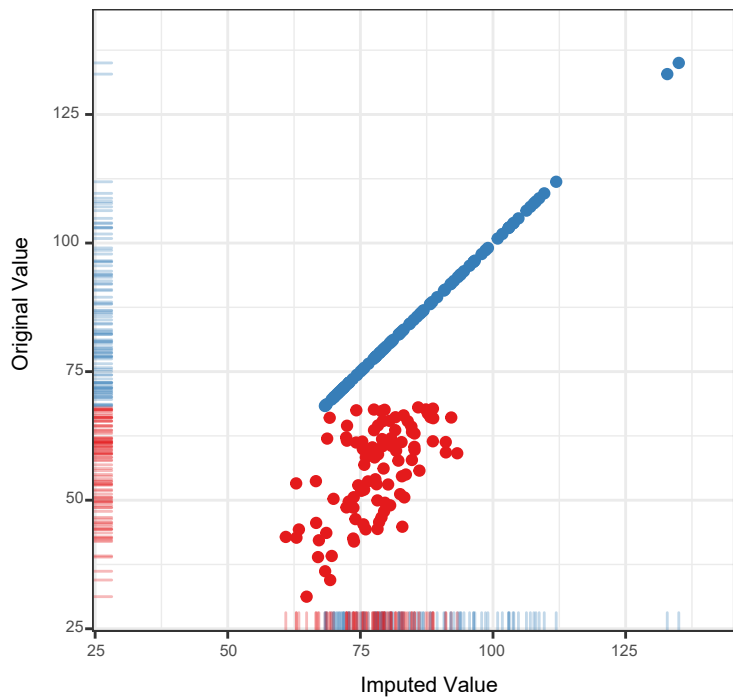
HM



QRILC



kNN-TN



GSimp

