

Comparative Assessment and Outlook on Methods for Imputing Proteomics Data

Minjie Shen

Virginia Polytechnic Institute and State University

Yi-Tan Chang

Virginia Polytechnic Institute and State University

Chiung-Ting Wu

Virginia Polytechnic Institute and State University

Sarah J Parker

Cedars Sinai Medical Center

Georgia Saylor

Wake Forest University

Yizhi Wang

Virginia Polytechnic Institute and State University

Guoqiang Yu

Virginia Polytechnic Institute and State University

Jennifer E. Van Eyk

Cedars Sinai Medical Center

Robert Clarke

University of Minnesota

David M. Herrington

Wake Forest University

Yue Wang (✉ yuewang@vt.edu)

Virginia Polytechnic Institute and State University

Research Article

Keywords: quantitative proteomics data analysis, evaluation methodologies, low-rank matrix factorization framework

Posted Date: March 18th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-298864/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Comparative assessment and outlook on methods for** 2 **imputing proteomics data**

3
4 Minjie Shen^{1,*}, Yi-Tan Chang^{1,*}, Chiung-Ting Wu^{1,*}, Sarah J Parker², Georgia
5 Saylor³, Yizhi Wang¹, Guoqiang Yu¹, Jennifer E. Van Eyk², Robert Clarke⁴,
6 David M. Herrington³, and Yue Wang^{1,†}

7 *Equal contribution; † Author for correspondence

8 ¹Department of Electrical and Computer Engineering, Virginia Polytechnic Institute and State
9 University, Arlington, VA 22203, USA; ²Advanced Clinical Biosystems Research Institute, Cedars
10 Sinai Medical Center, Los Angeles, CA 90048, USA; ³Department of Internal Medicine, Wake Forest
11 University, Winston-Salem, NC 27157, USA; ⁴The Hormel Institute, University of Minnesota, Austin,
12 MN 55912, USA

13

14 **Abstract**

15 **Background:** Missing values are a major issue in quantitative proteomics data
16 analysis. While many methods have been developed for imputing missing values in
17 high-throughput proteomics data, comparative assessment on the accuracy of
18 existing methods remains inconclusive, mainly because the true missing
19 mechanisms are complex and the existing evaluation methodologies are imperfect.
20 Moreover, few studies have provided an outlook of current and future development.

21 **Results:** We first report an assessment of eight representative methods collectively
22 targeting three typical missing mechanisms. The selected methods are compared on
23 both realistic simulation and real proteomics datasets, and the performance is
24 evaluated using three quantitative measures. We then discuss fused regularization
25 matrix factorization, a popular low-rank matrix factorization framework with similarity

26 and/or biological regularization, which is extendable to integrating multi-omics data
27 such as gene expressions or clinical variables. We further explore the potential
28 application of convex analysis of mixtures, a biologically-inspired latent variable
29 modeling strategy, to missing value imputation. The preliminary results on proteomics
30 data are provided together with an outlook into future development directions.

31 **Conclusion:** While a few winners emerged from our comparative assessment, data-
32 driven evaluation of imputation methods is imperfect because performance is
33 evaluated indirectly on artificial missing or masked values not authentic missing
34 values. Imputation accuracy may vary with signal intensity. Fused regularization
35 matrix factorization provides a possibility of incorporating external information.
36 Convex analysis of mixtures presents a biologically plausible new approach.

37

38 **Background**

39 Liquid chromatography coupled to mass spectrometry (LC-MS) is a popular method
40 for high-throughput identification and quantification of thousands of proteins in a single
41 analysis [1, 2]. The LC-MS signals can be displayed in a three-dimensional space with
42 the mass-to-charge ratios, retention times and intensities for the observed peptides.
43 However, this approach suffers from many missing values at the peptide or protein
44 level, which significantly reduces the amount of quantifiable proteins with an average
45 of 44% missing values [3-5].

46 While there are multiple causes for this missingness, three typical missing
47 mechanisms are widely acknowledged. Low abundant proteins may be missing because
48 their concentration is below the lower limit of detection (LLD); while poorly ionizing
49 peptides may cause proteins to be Missing Not at Random (MNAR) [6]. However,
50 missingness may also extend to mid- and even high-range intensities, statistically

51 categorized into Missing at Random (MAR) and Missing Completely at Random
52 (MCAR) [7]. MAR is actually missing conditionally at random given the observed,
53 known covariates, or even unknown covariates. MCAR depends neither on observed
54 nor on the missing data, thus the incomplete data are representative for the entire data.
55 Practically, MAR and MNAR cannot be distinguished because by definition missing
56 values are unknown [8]. More importantly, missing values in reality can originate from
57 a mix of both known and unknown missing mechanisms [7, 9].

58 A common solution for missingness is to impute the missing values based on
59 assumed missing mechanisms. But, this comes at the expense of potentially introducing
60 profound change in the distribution of protein-level intensities, because most of existing
61 methods are designed specifically for a single missing mechanism. This can have
62 unpredictable effects on downstream differential analyses. Moreover, while many
63 imputation methods have been adopted for imputing missing values in proteomics data,
64 comparative evaluation on their relative performance remains largely inconclusive, and
65 few studies provide an outlook addressing unresolved problems or future development
66 directions [4, 9, 10].

67 To gain first-hand insight into the strengths and limitations of both imputation
68 methods and assessment designs, we conduct a collective assessment of eight
69 representative methods involving three typical missing mechanisms in conjunction with
70 authentic missing values. Compared on a set of realistic and preserving simulations
71 derived from real proteomics data sets, the performance of the selected methods is
72 measured by three criteria, root-mean-square error (RMSE), normalized root-mean-
73 square error (NRMSE), and Sum of Ranks (SOR). There are several important
74 observations from this comparison study. First, while imputation methods perform
75 differentially under various missing mechanisms, algorithmic parameter settings, and

76 preprocessing procedures, there are a few methods that consistently outperformed peer
77 methods across a range of realistic simulation studies. Second, the quality of
78 performance assessment depends on the efficacy of simulation designs and a more
79 realistic simulation design should include authentic missing values and preserve
80 original overall data distribution. Third, existing assessment methodology is imperfect
81 in that performance is indirectly assessed on imputing either artificial or masked, but
82 not authentic missing values (see Discussion section).

83 To explore a more integrative strategy for improving imputation performance,
84 we discuss a low-rank matrix factorization framework with fused regularization on
85 sparsity and similarity – Fused Regularization Matrix factorization (FRMF) [11-13],
86 which can naturally integrate other-omics data such as gene expression or clinical
87 variables. We also introduce a biologically-inspired latent variable modeling strategy -
88 Convex analysis of Mixtures (CAM) [13, 14], which performs data imputation using
89 the original intensity data (before log-transformation). The preliminary results on real
90 proteomics data are provided together with an outlook into future development
91 directions.

92 **Results**

93 **Experimental design and protocol**

94 We selected eight representative methods for comparative assessment, based on their
95 intended missing mechanism(s) and imputation principles, summarized in **Figure 1**.
96 One method (Min/2) is devoted to MNAR (LLD) [7], two methods (swKNN and
97 pwKNN) are tailored to MAR (local-similarity) [15], and five methods (Mean, PPCA,
98 NIPALS, SVD, and SVT) are designed for MCAR/MAR (global-structure or low-rank
99 matrix factorization) [7, 9, 16-18]. We then explored and tested several variants of

100 FRMF and CAM, where local similarity information is obtained from baseline or other
101 data acquired from the same samples.

102 We conducted the comparative assessments in two complementary simulation
103 settings. First, the realistic simulation data were generated from the observed data
104 portion (no authentic missing value) of a real proteomics dataset, where artificial
105 missing values were introduced involving two typical missing mechanisms and used
106 for performance assessment. Second, the realistic simulation data were generated from
107 the complete data matrix (including authentic missing values) of a real proteomics
108 dataset, where a small percentage of data points were randomly set-aside (masked
109 values) and used solely for performance assessment. The preprocessing eliminates
110 those proteins whose missing rates are higher than 80% and then performs log₂
111 transformation [19]. The parameters were optimized for each imputation method by
112 parameter sweeping over a wide range of settings at each missing rate. The overall
113 experimental workflow is given in **Figure 2**.

114 **Real proteomics data**

115 The real LC-MS proteomics data form the base from which the simulation data sets
116 were produced [6]. The data were acquired using data-independent acquisition (DIA)
117 protocol, and protein level output was generated by mapDIA [20]. The dataset contains
118 200 samples associated with 2,682 proteins measured in human left anterior descending
119 (LAD) coronary arteries collected as part of a study of coronary and aortic
120 atherosclerosis [21]. The data were produced in three separate batches, indexed by A,
121 B, and C, and all have passed quality control and preprocessing procedures,
122 summarized in **Table 1** (Supplementary Information).

123 **Table 1.** Summary of real proteomics datasets used in this work.

	Sample size	Protein size	Total Missing Rate #MV/(#Sample*#Protein)	Setting #1 protein size (non-missing proteins)	Setting #2 protein size (proteins with <= 80% missing rate)
Batch A	98	2107	24.67%	751 (35.64%)	1935 (91.84%)
Batch B	55	2604	29.63%	819 (31.45%)	2324 (89.25%)
Batch C	47	2590	25.52%	976 (37.68%)	2325 (89.77%)

124

125 **Simulation data generated from the observed portion of data matrix**

126 Based on the observed data portion (no authentic missing values), we adopted a hybrid
 127 missing data model and used the R package `imputeLCMD` to introduce artificial
 128 missing values while preserving the original observed data patterns [22]. Specifically,
 129 MCAR missing values were introduced by randomly replacing some data points with
 130 ‘NA’ (not available) according to the designed missing rates (approximately from 1%
 131 to 50%); MNAR missing values were introduced by quantile cut-off for the full data
 132 set [7, 9, 19]; and mixed MCAR and MNAR missing values were introduced by
 133 assigning $(1 - \beta)$ portion of MCAR and β portion of MNAR; corresponding to
 134 missing rate α and $\beta = 0, 0.1, 1$ (Supplementary Information).

135 **Simulation data with set-aside masked values from the full data matrix**

136 In this simulation setting, we used the full data matrix (including both observed and
 137 authentic missing values) from the human coronary proteomics data set. To preserve
 138 the original patterns of both observed and authentic missing data, for each protein, a
 139 small percentage of data points in the complete data matrix were randomly set-aside as
 140 ‘NA’ (masked values) with the masking rate(s) proportional to the authentic missing
 141 rate(s). This procedure was repeated for all proteins and the masked values were

142 considered as a mix of MNAR and MAR conditioned on the observed missing rates
143 and data patterns (**Figure 3**, Supplementary Information).

144 **Performance assessment focused on MNAR**

145 As shown in **Figure 4** (see additional results in Supplementary Information), under
146 MNAR missing mechanism assumption, SVT and Min/2 yielded the best performance
147 in both simulation settings, the relative performance of SVT and Min/2 depends on the
148 missing rates and criterion used for evaluation. The MNAR-devoted method, Min/2,
149 performs much better than the others as expected while the baseline method, Mean,
150 performance is the worst among all methods (see additional results in Supplementary
151 Information). Note that SOR increases expectedly when the missing rate increases
152 because SOR is positively associated with the number of missing proteins, therefore the
153 value of SOR may not imply the absolute performance of a method.

154 **Performance assessment focused on MCAR**

155 The imputation performance of the eight methods on MCAR mechanism is shown in
156 **Figure 5** (see additional results in Supplementary Information). The experimental
157 results show that NIPALS outperforms all other methods in both simulation settings
158 and for almost all three evaluation criteria; while SVT is the best when RMSE is used;
159 Min/2 performs the worst among all other methods in all cases that may be expected
160 due to its design for MNAR mechanism; and Mean performs flatly over different
161 missing rates with an expected baseline performance except for Min/2. In addition,
162 while all methods perform worse when total missing rate increases, the ranking of their
163 relative performances remains unchanged (see additional results in Supplementary
164 Information).

165 Performance assessment focused on authentic missing values

166 As shown in **Figure 6** (see additional results in Supplementary Information), NIPALS,
167 SVT, and protein-wise or sample-wise KNN achieve the best performance, where
168 authentic missing values are dominant and imputation accuracy is evaluated on the
169 masked values. MNAR-devoted method Min/2 performs the worst as expected. Similar
170 to the case of MCAR, low-rank methods and local-similarity methods perform worse
171 when total missing rate increases, and among these methods, SVD and PPCA perform
172 even worse than the baseline method Mean when total missing rate is large. Note that
173 because low abundant proteins often have higher authentic missing rates and
174 accordingly higher masking rates, more low abundant proteins (possibly the minimum
175 values) are masked over highly expressed proteins, and the counterintuitive decrease in
176 NRMSE by Min/2 is expected when the authentic missing rate increases.

177 Evaluation of the FRMF method focused on authentic missing values

178 We evaluated three variants of the FRMF method. RMF serves as a baseline regularized
179 matrix factorization algorithm; FRMF_self introduces a fused-regularization utilizing
180 the similarity among samples embedded within data matrix; and FRMF_cross_patho
181 exploits external pathological scores again via fused-regularization strategy where the
182 pathological scores are the qualitative percentages of the intimal surface involvement
183 of various atherosclerotic changes graded by pathologists [6].

184 The experimental results are shown in **Figure 7**. While RMF performs
185 comparably and expectedly to SVD, both FRMF_self and FRMF_cross_patho
186 significantly outperform RMF. This preliminary result indicates a potential benefit of
187 combining global low-rank and local-similarly regularizations, as well as leveraging
188 external information via fused regularization.

189 **Evaluation of the CAM method focused on authentic missing values**

190 Based on biologically-inspired latent variable modeling of complex tissues - CAM [13,
191 14], we proposed and evaluated three variants of the CAM based imputation method.
192 CAM_complete performs CAM based imputation using the non-missing portion of full
193 data matrix; CAM_SVT and CAM_NIPALS perform CAM based imputation using full
194 data matrix while initialized by SVT and NIPALS, respectively.

195 The experimental results are shown in **Figure 8**. Expectedly, CAM_complete
196 performs much better than the baseline method Mean. More importantly, both
197 CAM_NIPALS and CAM_SVT consistently outperform NIPALS and SVT - the two
198 top performers indicated in our earlier comparative assessment. This preliminary result
199 suggests that biologically-plausible latent variable modeling may potentially improve
200 imputation accuracy within the framework of low-rank optimization.

201 **Method**

202 **Brief introduction to the eight existing methods**

- 203 • **Min/2 (half minimum):** Assuming the MNAR missing mechanism, for each
204 protein, replacing missing values with half the minimum value of observed
205 intensities in that protein across samples [6, 10].
- 206 • **Mean:** Assuming the MAR/MCAR missing mechanism, for each protein,
207 replacing missing values with mean value of observed intensities in that protein
208 across samples [6, 10].
- 209 • **swKNN (sample-wise k-nearest neighbors):** Assuming the MAR missing
210 mechanism and leveraging local similarity among samples, for each protein,
211 replacing missing values with weighted average of observed intensities in that
212 protein proportional to the proximities of k-nearest neighboring samples [10].

- 213 • **pwKNN (protein-wise k-nearest neighbors):** Assuming the MAR missing
214 mechanism and leveraging local similarity among proteins, for each sample,
215 replacing the missing values with weighted average of observed intensities in
216 that sample proportional to the proximities of k-nearest neighboring proteins
217 (with protein-wise normalization) [10].
- 218 • **PPCA (probabilistic PCA):** Assuming the MCAR/MAR missing mechanism,
219 a low-rank probabilistic PCA matrix factorization is estimated by the
220 expectation maximization (EM) algorithm and subsequently used to impute
221 missing values [23].
- 222 • **NIPALS (non-linear estimation by iterative partial least squares):**
223 Assuming the MCAR/MAR missing mechanism, a low-rank missing-data-
224 tolerant PCA matrix factorization is estimated by iterative regression and
225 subsequently used to impute missing values [24, 25].
- 226 • **SVD (SVDImpute):** Assuming the MCAR/MAR missing mechanism, a low-
227 rank SVD matrix factorization is estimated by the EM algorithm and
228 subsequently used to impute missing values [24, 26].
- 229 • **SVT (singular value thresholding):** Assuming MCAR/MAR missing
230 mechanism, a low-rank SVT matrix factorization is estimated by iteratively
231 solving a nuclear norm minimization problem and subsequently used to impute
232 missing values [18].

233 **Performance measures**

234 Three quantitative measures are used to evaluate imputation accuracy, namely
235 Root Mean Square Error (RMSE), Normalized Root Mean Square Error (NRMSE), and
236 Sum of Ranks (SOR). Specifically, RMSE and NRMSE are given by [27, 28]

$$237 \quad \text{RMSE} = \sqrt{\frac{\sum_{\Omega} (\hat{X}_{\Omega} - X_{\Omega})^2}{|\Omega|}}, \quad \text{NRMSE} = \sqrt{\frac{\sum_{\Omega} (\hat{X}_{\Omega} - X_{\Omega})^2}{|\Omega| \sigma_{X_{\Omega}}^2}},$$

238 respectively, where Ω is the index set of missing values in complete data matrix X , $|\Omega|$
 239 is the total number of missing values, \hat{X} is the imputed complete data matrix, and $\sigma_{X_{\Omega}}^2$
 240 is the variance of missing values. To address the bias of NRMSE under MNAR missing
 241 mechanism, SOR has been proposed as [19]

$$242 \quad \text{SOR} = \sum_{i=1}^P \text{rank}(\text{NRMSE}_i),$$

243 where P is the number of proteins containing at least one missing value, i is the protein
 244 index in this protein subset, and $\text{rank}(\text{NRMSE}_i)$ is the ranks of protein-wise NRMSE
 245 across different imputation methods.

246 Introduction to FRMF method

247 As aforementioned, low-rank matrix factorization has been a popular and effective
 248 approach for missing data imputation [12]. For imputing proteomics data, the
 249 assumption is that there is only a small number of biological processes determining the
 250 expression profiles. Consider an $m \times n$ complete data matrix X describing m samples
 251 and n proteins. A low-rank matrix factorization approach seeks to approximate X
 252 containing missing values by a linear latent variable model,

$$253 \quad X_{m \times n} = A_{m \times l} \times S_{l \times n}, \quad (1)$$

254 where $A_{m \times l}$ and $S_{l \times n}$ are the low-rank factor matrices, and $l \ll \min(m, n)$. In order to
 255 prevent overfitting, the solution is often formulated as a regularized sparse SVD
 256 minimization problem on the observed values

$$257 \quad \min \sum_{i=1}^m \sum_{j=1}^n I(X_{ij} \neq \text{NA}) (X_{ij} - A_i S_j)^2 + \lambda_A \|A\|_F^2 + \lambda_S \|S\|_F^2,$$

258 where $\|\cdot\|_F^2$ denotes the Frobenius norm, $I(\cdot)$ is the indicator function, and $\lambda_A, \lambda_S > 0$
 259 are the regularization parameters. When local similarity information is available, FRMF
 260 can be formulated by adding a fused regularization term

$$261 \quad \min \sum_{i=1}^m \sum_{j=1}^n I(X_{ij} \neq \text{NA})(X_{ij} - A_i S_j)^2 + \lambda_A \|A\|_F^2 + \lambda_S \|S\|_F^2$$

$$262 \quad + \alpha \sum_{i=1}^m \sum_{k \in \mathcal{F}(i)} \|A_i - A_k\|_F^2,$$

263 where α is the fused regularization parameter, and $\mathcal{F}(i)$ denotes the neighborhood
 264 sample subset of sample i and can be determined using baseline data or other relevant
 265 measurements e.g. gene expression or pathological score. In our study, $\mathcal{F}(i)$ is
 266 determined by the between-sample cosine similarity $\cos(X_i, X_k)$ based on data matrix
 267 in FRMF_self, or $\cos(P_i, P_k)$ based on pathological scores in FRMF_cross_patho.

268 Introduction to CAM method

269 CAM is a latent variable modeling and deconvolution technique previously used for
 270 identifying biologically-interpretable cell subtypes $S_{l \times n}$ and their composition $A_{m \times l}$ in
 271 complex tissues [6, 13, 14, 21]. We adopt the CAM framework into (1) and demonstrate
 272 that hybrid CAM_SVT and CAM_NIPALS can handle missing values naturally and
 273 this combination leads to a novel and biologically-plausible imputation strategy. The
 274 workflow of CAM based method with three variants is given in **Figure 9**.

275 Discussion

276 The quality of simulating assumed missing mechanisms (MNAR and MCAR) depends
 277 on the efficacy of simulation tools. However, because the simulation uses only the
 278 observed portion of data matrix, the introduced artificial missing values cannot fully
 279 resemble authentic missing mechanisms and/or patterns in relation to original overall

280 data distribution. More critically, performance is actually assessed on imputing
281 artificial not authentic missing values, where the overall data distribution may be
282 distorted.

283 To address the aforementioned issues in the presence of authentic missing
284 values, a small percentage of set-aside values are introduced into complete data matrix
285 and used solely for assessment purpose. Because masked values are randomly assigned
286 onto both observed and authentic missing values, the simulation maximally preserves
287 original overall data distribution. Note that masked values may represent a mix of
288 MNAR (high missing rate associated with low protein abundance) and MAR (joint
289 distribution of both observed and authentic missing values). However, performance is
290 assessed indirectly on imputing masked not authentic missing values.

291 Imputation accuracy would be arguably affected by data preprocessing and
292 algorithmic parameter setting. In this study, sample-wise normalization and protein-
293 wise standardization are performed based on the requirements of each method. While
294 these preprocessing notably affects the scale of NRMSE, relative performances across
295 methods remain consistent. The experimental results show that imputation performance
296 varies with parameter setting, while there appears no theoretical guideline for
297 optimizing parameter setting.

298 While FRMF is a promising and novel imputation approach, its effectiveness
299 for improving classic low-rank methods would depend on diversity among samples,
300 discriminatory power of similarity measure, and complementary nature of additional
301 and relevant measurements. Newly proposed CAM method represents an interesting
302 direction for further development. More importantly, CAM performs missing value
303 imputation using original intensity rather than log-transformed data, and this is

304 mathematically more rigorous because log-transformation violates the linear nature of
305 low-rank matrix factorization [29].

306

307 **Declarations**

308 - Ethics approval and consent to participate

309 Not applicable

310 - Consent for publication

311 Not applicable

312 - Availability of data and materials

313 The scripts used in the paper is available in R script ProImput.

314 Code for all experiments can be found in the vignette at

315 <https://github.com/MinjieSh/ProImput>. The operation system

316 can be any system supporting R language.

317 - Competing interests

318 The authors declare that they have no competing interests.

319 - Funding

320 This work has been supported by the National Institutes of

321 Health under Grants HL111362-05A1, HL133932, NS115658-

322 01, and the Department of Defense under Grant W81XWH-18-

323 1-0723 (BC171885P1).

324 - Authors' Contributions

325 MJS, YTC, CTW and YW initiated this research and developed

326 the methods. MJS and YW drafted the manuscript. MJS and

327 CTW developed the assessment procedure and software. SJP

328 and GS contributed with biomedical case study. YZW, GQY,
 329 JEVE and DMH contributed with biomedical significance and
 330 discussion of the work. All authors have read, commented on
 331 and accepted the final manuscript.

332 - Acknowledgements

333 Not applicable

334 References

- 335 1. Canterbury JD, Merrihew GE, MacCoss MJ, Goodlett DR, Shaffer SA:
 336 **Comparison of data acquisition strategies on quadrupole ion trap**
 337 **instrumentation for shotgun proteomics.** *Journal of The American Society*
 338 *for Mass Spectrometry* 2014, **25**(12):2048-2059.
- 339 2. Doerr A: **DIA mass spectrometry.** *Nature methods* 2014, **12**(1):35.
- 340 3. Goeminne LJE, Sticker A, Martens L, Gevaert K, Clement L: **MSqRob Takes**
 341 **the Missing Hurdle: Uniting Intensity- and Count-Based Proteomics.** *Anal*
 342 *Chem* 2020, **92**(9):6278-6287.
- 343 4. Ma W, al. e, Wang P: **DreamAI: algorithm for the imputation of**
 344 **proteomics data.** *bioRxiv* 2020.
- 345 5. Dabke K, Kreimer S, Jones MR, Parker SJ: **A Simple Optimization**
 346 **Workflow to Enable Precise and Accurate Imputation of Missing Values**
 347 **in Proteomic Datasets.** *bioRxiv* 2020.
- 348 6. Herrington DM, Mao C, Parker SJ, Fu Z, Yu G, Chen L, Venkatraman V, Fu
 349 Y, Wang Y, Howard TD *et al*: **Proteomic Architecture of Human Coronary**
 350 **and Aortic Atherosclerosis.** *Circulation* 2018, **137**(25):2741-2756.
- 351 7. Lazar C, Gatto L, Ferro M, Bruley C, Burger T: **Accounting for the multiple**
 352 **natures of missing values in label-free quantitative proteomics data sets to**
 353 **compare imputation strategies.** *Journal of proteome research* 2016,
 354 **15**(4):1116-1125.
- 355 8. Jakobsen JC, Gluud C, Wetterslev J, Winkel P: **When and how should**
 356 **multiple imputation be used for handling missing data in randomised**
 357 **clinical trials - a practical guide with flowcharts.** *BMC Med Res Methodol*
 358 2017, **17**(1):162.
- 359 9. Webb-Robertson B-JM, Wiberg HK, Matzke MM, Brown JN, Wang J,
 360 McDermott JE, Smith RD, Rodland KD, Metz TO, Pounds JG: **Review,**
 361 **evaluation, and discussion of the challenges of missing value imputation**
 362 **for mass spectrometry-based label-free global proteomics.** *Journal of*
 363 *proteome research* 2015, **14**(5):1993-2001.
- 364 10. Liu M, Dongre A: **Proper imputation of missing values in proteomics**
 365 **datasets for differential expression analysis.** *Brief Bioinform* 2020.
- 366 11. Lin X, Boutros PC: **Optimization and expansion of non-negative matrix**
 367 **factorization.** *BMC bioinformatics* 2020, **21**(1):7.

- 368 12. Ma H, Zhou D, Liu C, Lyu MR, King I: **Recommender systems with social**
369 **regularization**. In: *The fourth ACM international conference on Web search*
370 *and data mining: 2011; Hong Kong*. ACM Press: 287-296.
- 371 13. Wang N, Hoffman EP, Chen L, Chen L, Zhang Z, Liu C, Yu G, Herrington
372 DM, Clarke R, Wang Y: **Mathematical modelling of transcriptional**
373 **heterogeneity identifies novel markers and subpopulations in complex**
374 **tissues**. *Scientific Reports* 2016, **6**:18909.
- 375 14. Chen L, Wu CT, Wang N, Herrington DM, Clarke R, Wang Y: **debCAM: a**
376 **bioconductor R package for fully unsupervised deconvolution of complex**
377 **tissues**. *Bioinformatics (Oxford, England)* 2020, **36**(12):3927-3929.
- 378 15. Rahman SA, Huang Y, Claassen J, Heintzman N, Kleinberg S: **Combining**
379 **Fourier and lagged k-nearest neighbor imputation for biomedical time**
380 **series data**. *Journal of biomedical informatics* 2015, **58**:198-207.
- 381 16. Pedersen AB, Mikkelsen EM, Cronin-Fenton D, Kristensen NR, Pham TM,
382 Pedersen L, Petersen I: **Missing data and multiple imputation in clinical**
383 **epidemiological research**. *Clinical epidemiology* 2017, **9**:157.
- 384 17. John C, Ekpenyong EJ, Nworu CC: **Imputation of missing values in**
385 **economic and financial time series data using five principal component**
386 **analysis approaches**. *CBN Journal of Applied Statistics* 2019, **10**(1):51-73.
- 387 18. Cai J-F, Candès EJ, Shen Z: **A singular value thresholding algorithm for**
388 **matrix completion**. *SIAM Journal on optimization* 2010, **20**(4):1956-1982.
- 389 19. Wei R, Wang J, Su M, Jia E, Chen S, Chen T, Ni Y: **Missing value**
390 **imputation approach for mass spectrometry-based metabolomics data**.
391 *Scientific reports* 2018, **8**(1):1-10.
- 392 20. Teo G, Kim S, Tsou C-C, Collins B, Gingras A-C, Nesvizhskii AI, Choi H:
393 **mapDIA: Preprocessing and statistical analysis of quantitative proteomics**
394 **data from data independent acquisition mass spectrometry**. *Journal of*
395 *proteomics* 2015, **129**:108-120.
- 396 21. Parker SJ, Chen L, Spivia W, Saylor G, Mao C, Venkatraman V, Holewinski
397 RJ, Mastali M, Pandey R, Athas G *et al*: **Identification of Putative Early**
398 **Atherosclerosis Biomarkers by Unsupervised Deconvolution of**
399 **Heterogeneous Vascular Proteomes**. *J Proteome Res* 2020, **19**(7):2794-
400 2806.
- 401 22. **imputeLCMD: A collection of methods for left-censored missing data**
402 **imputation** [<https://cran.r-project.org/package=imputeLCMD>]
- 403 23. Tipping ME, Bishop CM: **Probabilistic principal component analysis**.
404 *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*
405 1999, **61**(3):611-622.
- 406 24. Stacklies W, Redestig H, Scholz M, Walther D, Selbig J: **pcaMethods—a**
407 **bioconductor package providing PCA methods for incomplete data**.
408 *Bioinformatics* 2007, **23**(9):1164-1167.
- 409 25. Ochoa-Muñoz AF, González-Rojas VM, Pardo CE: **Missing data in multiple**
410 **correspondence analysis under the available data principle of the**
411 **NIPALS algorithm**. *DYNA* 2019, **86**(211):249-257.
- 412 26. Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R,
413 Botstein D, Altman RB: **Missing value estimation methods for DNA**
414 **microarrays**. *Bioinformatics* 2001, **17**(6):520-525.
- 415 27. Stekhoven DJ, Bühlmann P: **MissForest—non-parametric missing value**
416 **imputation for mixed-type data**. *Bioinformatics* 2012, **28**(1):112-118.

- 417 28. Oba S, Sato M, Takemasa I, Monden M, Matsubara K, Ishii S: **A Bayesian**
418 **missing value estimation method for gene expression profile data.**
419 *Bioinformatics* 2003, **19**(16):2088-2096.
- 420 29. Zhong Y, Liu Z: **Gene expression deconvolution in linear space.** *Nat*
421 *Methods* 2011, **9**(1):8-9; author reply 9.
422

Figures

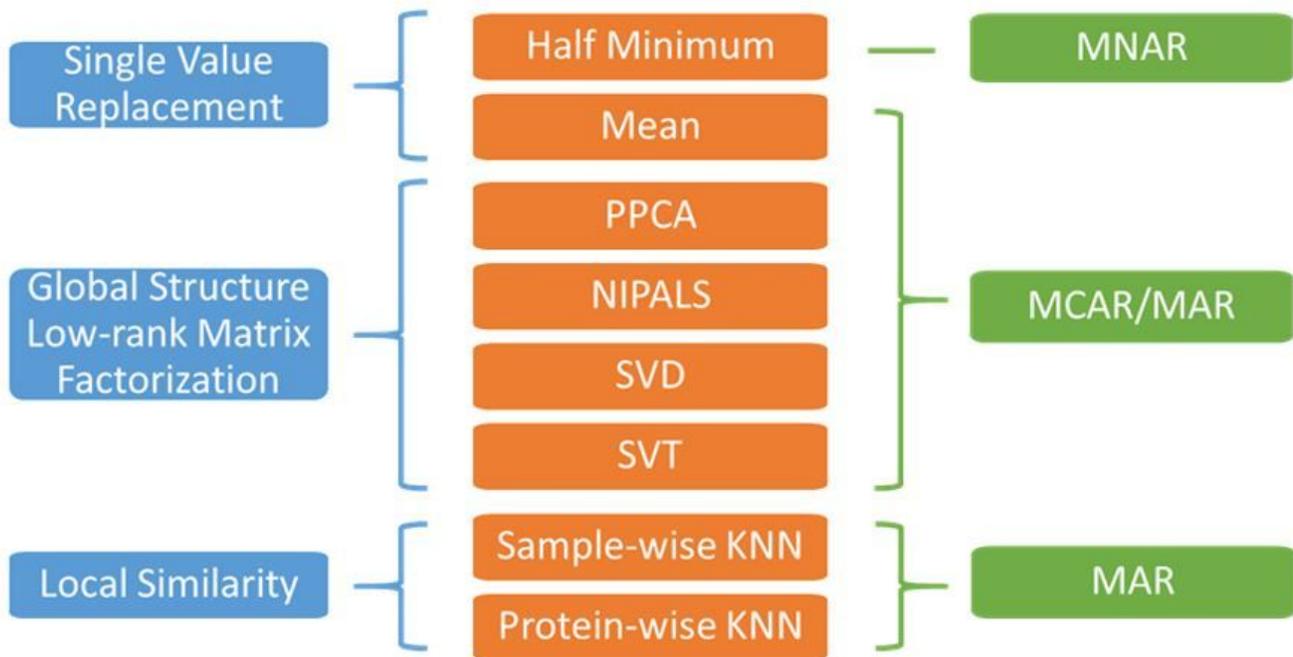


Figure 1

Comparative assessment of eight representative missing value imputation methods, divided into three categories.

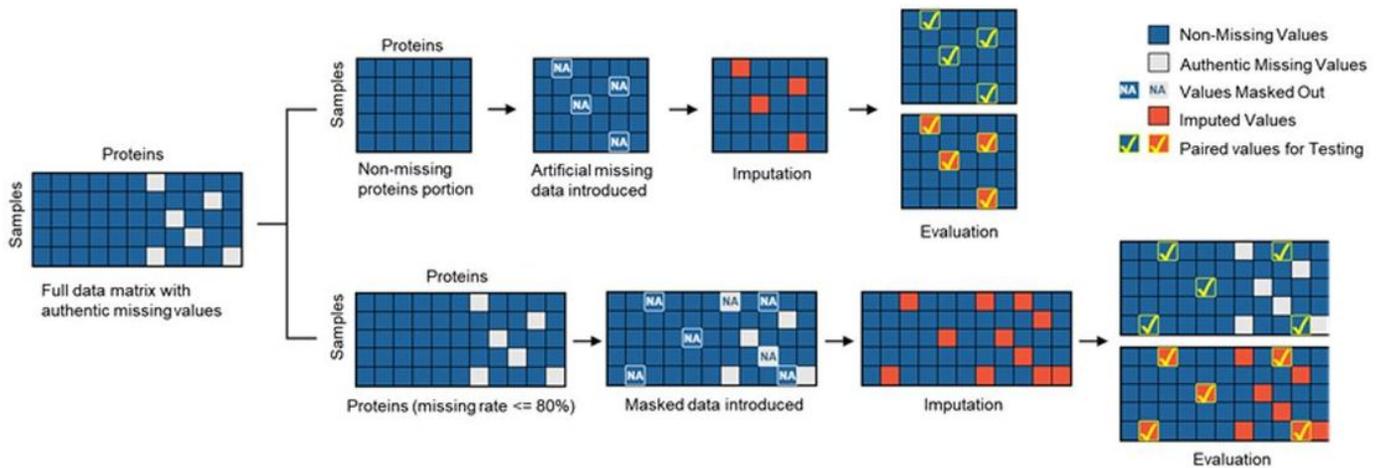


Figure 2

Two-phased workflow of realistic simulation-based assessment on missing value imputation methods.

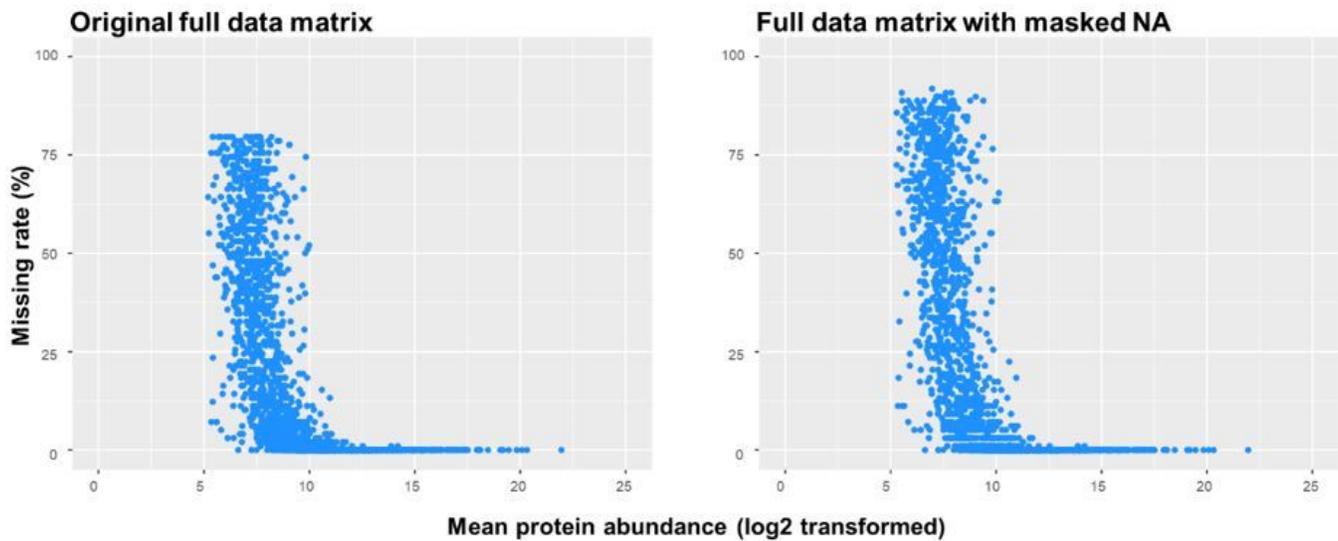


Figure 3

The overall pattern of missing values illustrated by the relationship between protein missing rate and protein mean intensity, before (left panel) and after (right panel) introducing masked NA.

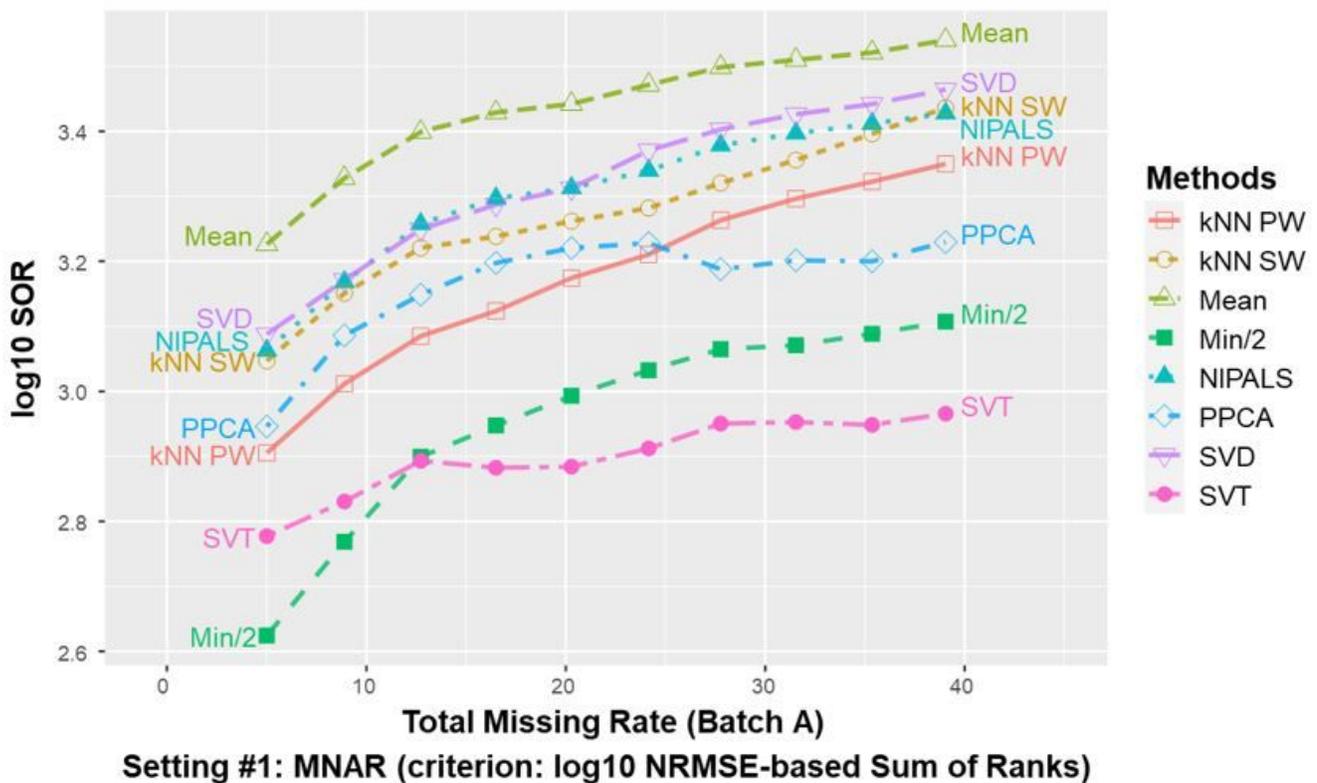


Figure 4

Imputation performance of the eight methods on the simulation data of setting #1, with assumed MNAR missing mechanism and varying total missing rates.

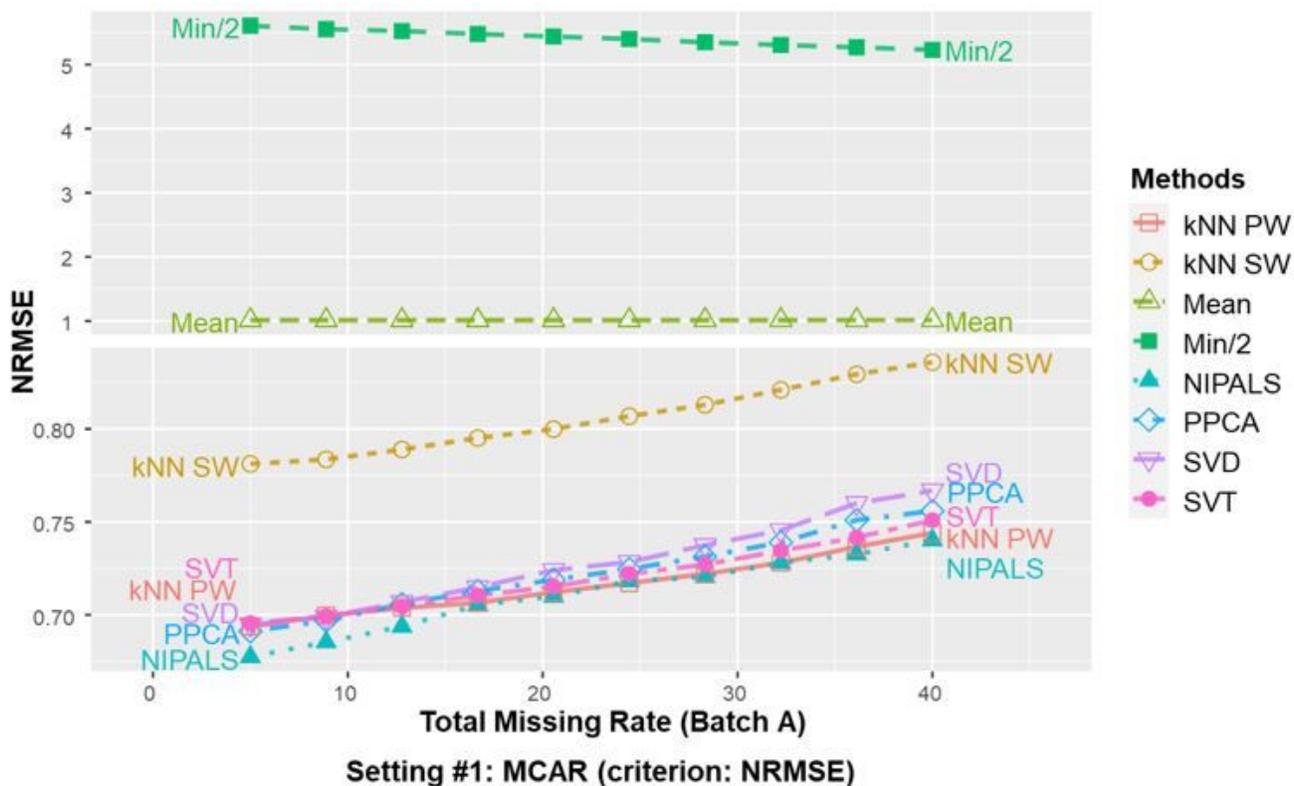


Figure 5

Imputation performance of the eight methods on the simulation data of setting #1, with assumed MCAR missing mechanism and varying total missing rates.

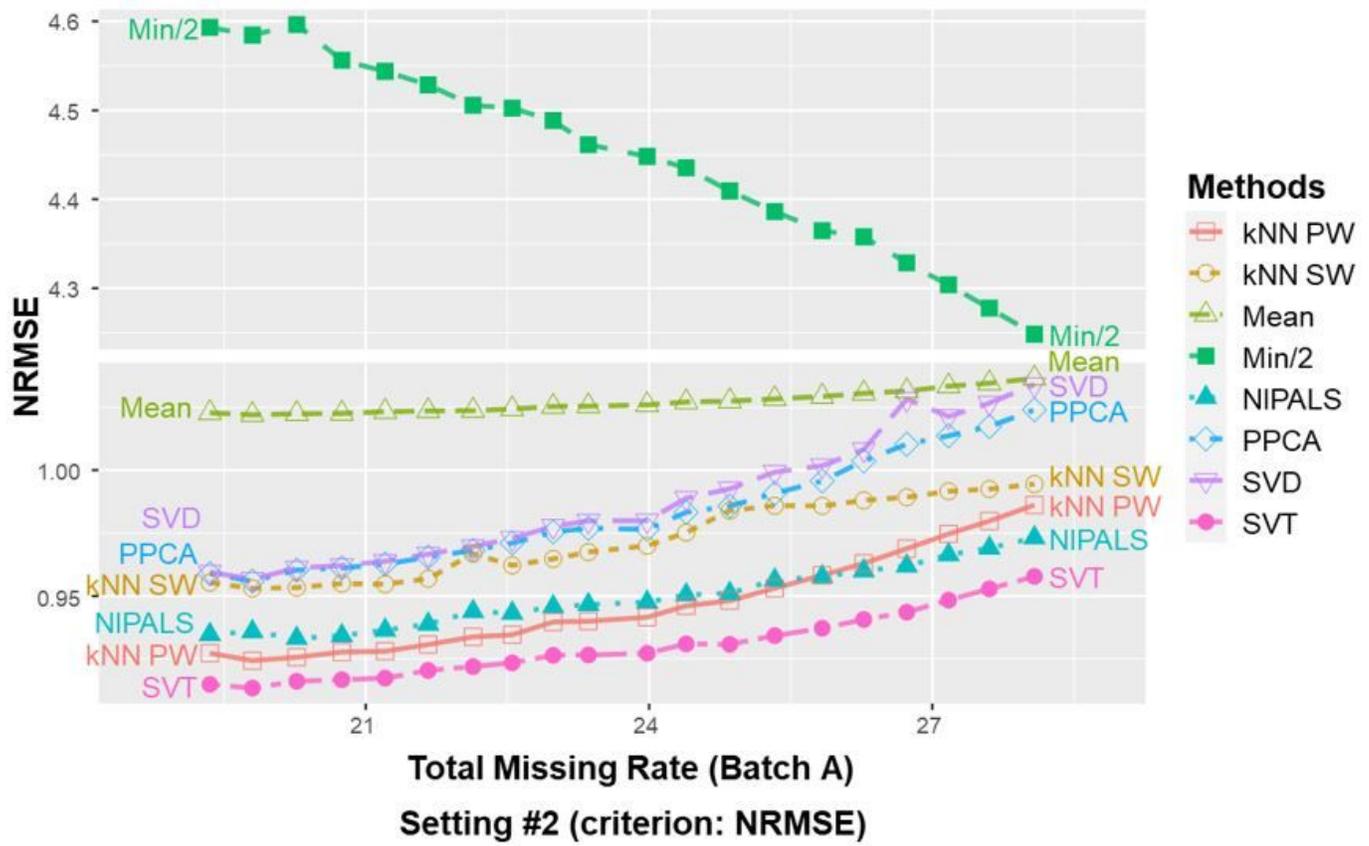


Figure 6

Imputation performance of the eight methods on the simulation data of setting #2, focusing on authentic missing mechanism and varying masked rates.



Figure 7

Imputation performance of the FRMF variants on the simulation data of setting #2, with varying masked rates.

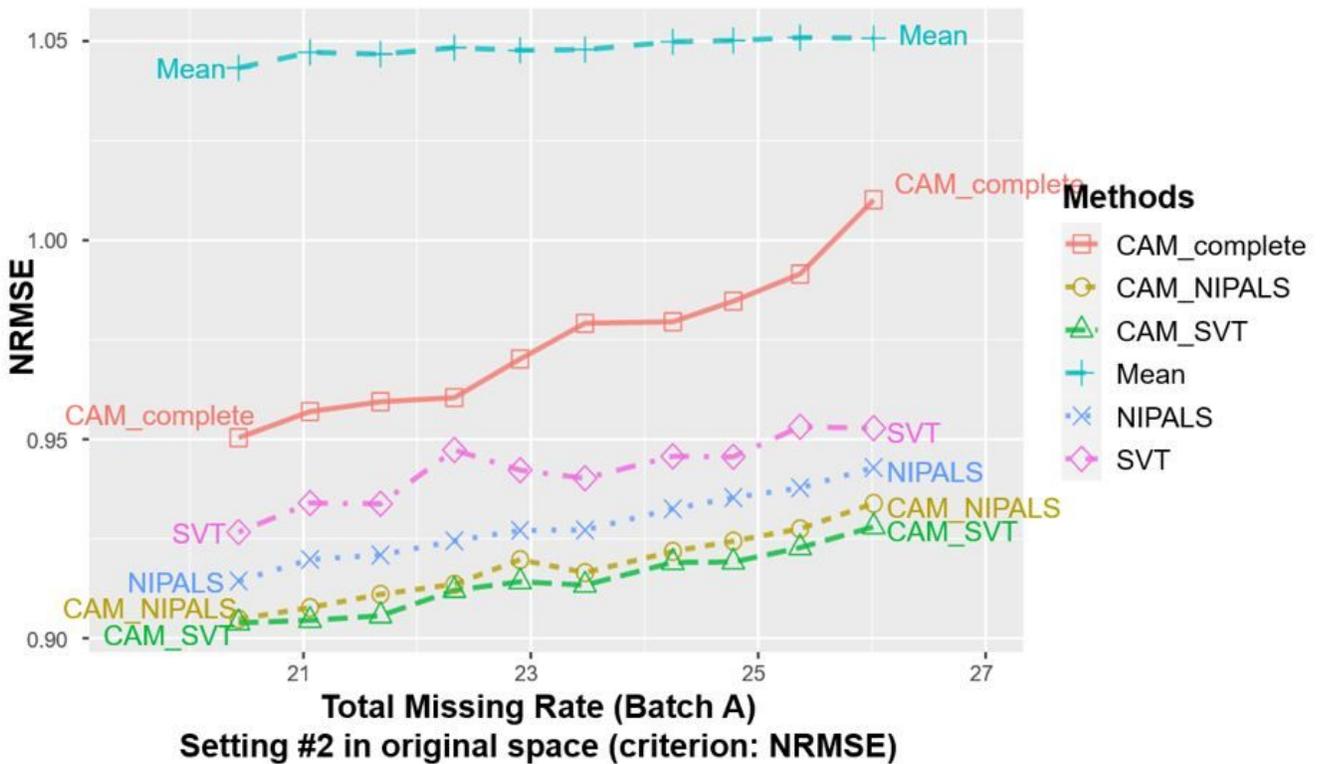


Figure 8

Imputation performance of CAM variants on the simulation data of setting #2, with varying masked rates, in comparison to that of Mean, SVT, and NIPALS. The imputation accuracy is evaluated in the original intensity space (before log-transformation).

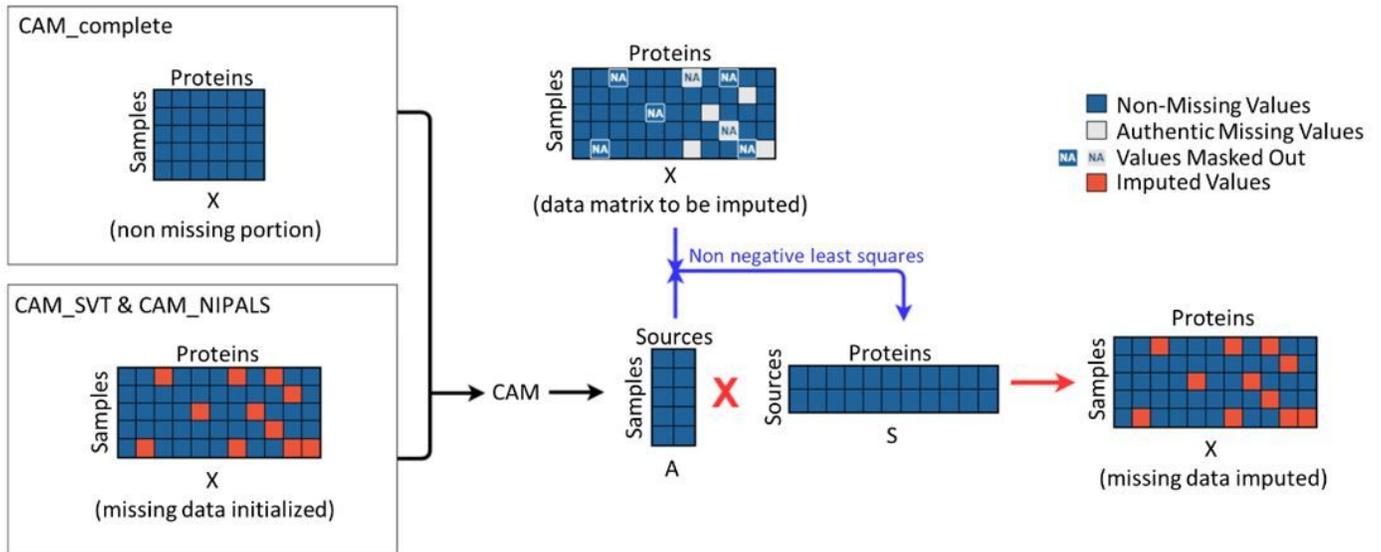


Figure 9

Workflow of the CAM based imputation method with two variant algorithms.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryInformation.docx](#)