

# Accounting for the Multiple Natures of Missing Values in Label-Free Quantitative Proteomics Data Sets to Compare Imputation Strategies

Cosmin Lazar,<sup>†,§,||</sup> Laurent Gatto,<sup>⊥,#</sup> Myriam Ferro,<sup>†,§,||</sup> Christophe Bruley,<sup>†,§,||</sup> and Thomas Burger<sup>\*,†,‡,§,||</sup>

<sup>†</sup>Univ. Grenoble Alpes, iRTSV-BGE, F-38000 Grenoble, France

<sup>‡</sup>CNRS, iRTSV-BGE, F-38000 Grenoble, France

<sup>§</sup>CEA, iRTSV-BGE, F-38000 Grenoble, France

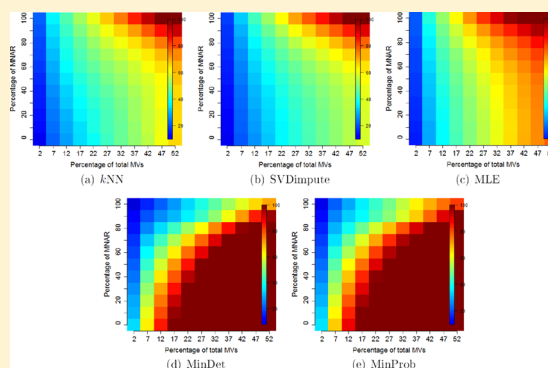
<sup>||</sup>INSERM, BGE, F-38000 Grenoble, France

<sup>⊥</sup>Computational Proteomics Unit, Cambridge CB2 1GA, United Kingdom

<sup>#</sup>Cambridge Center for Proteomics, Cambridge CB2 1GA, United Kingdom

**ABSTRACT:** Missing values are a genuine issue in label-free quantitative proteomics. Recent works have surveyed the different statistical methods to conduct imputation and have compared them on real or simulated data sets and recommended a list of missing value imputation methods for proteomics application. Although insightful, these comparisons do not account for two important facts: (i) depending on the proteomics data set, the missingness mechanism may be of different natures and (ii) each imputation method is devoted to a specific type of missingness mechanism. As a result, we believe that the question at stake is not to find the most accurate imputation method in general but instead the most appropriate one. We describe a series of comparisons that support our views: For instance, we show that a supposedly “under-performing” method (i.e., giving baseline average results), if applied at the “appropriate” time in the data-processing pipeline (before or after peptide aggregation) on a data set with the “appropriate” nature of missing values, can outperform a blindly applied, supposedly “better-performing” method (i.e., the reference method from the state-of-the-art). This leads us to formulate few practical guidelines regarding the choice and the application of an imputation method in a proteomics context.

**KEYWORDS:** label-free relative quantitative proteomics, missing value imputation



## 1. INTRODUCTION

The high rate of missing values in label-free quantitative proteomics is a major concern.<sup>1</sup> From the literature, in the case of LC–MS/MS approaches, it frequently ranges between 10 and 50%, while the proportion of peptides/proteins that exhibit at least one missing value can be very high, ranging in between 70 and 90%.<sup>2</sup> As a consequence, it was originally proposed to apply imputation methods originally developed for transcriptomics and microarray data analysis<sup>3</sup> to proteomics data. Then, more general methods, developed in a theoretical statistical context, were considered<sup>4</sup> and adapted to some extent to proteomics data sets.<sup>5</sup> To date, numerous methods exist and are available to any practitioner, either as independent packages<sup>6–8</sup> or through dedicated pipeline packages such as MSnbase.<sup>9</sup> In addition, several methods have been reported that successfully leverage on a multiomics context to impute proteomics missing values on the basis of transcriptomics observed values.<sup>10–12</sup> Recently, a comprehensive survey<sup>13</sup> compared and discussed some well-known imputation algorithms in the context of proteomics applications. There are numerous conclusions that can be drawn from this survey or from references therein.

First, there are multiple reasons why values are missing, accounting for biochemical and analytical (miscleavage, dynamic range, ionization competition, ion suppression, etc.) to bioinformatics mechanisms (peptide misidentification, ambiguous matching of the precursors in the quantitation step, etc.); however, regardless of their origins, missing values can be cast in three categories with regards to the statistical mechanisms that best describe them. In fact, statisticians have defined three types of missing values:<sup>4</sup>

- Missing Completely At Random (MCAR), which in a proteomics data set, correspond to the combination and propagation of multiple minor errors or stochastic fluctuations. (For instance, a misidentified peptide can or cannot be balanced by the alignment of the precursor maps, leading to an abundance value or, on the contrary, to a missing value). As a result, each missing value cannot be directly explained by the nature of the peptide or by its measured intensity.<sup>5</sup> As a result, MCAR affects the entire data set with a uniform distribution.

**Received:** October 21, 2015

**Published:** February 23, 2016

•Missing At Random (MAR), which is a more general class than MCAR, where conditional dependencies are accounted for. In a proteomics data set, it is classically assumed that all MAR values are also MCAR so that one is little interested in MAR;<sup>5</sup> however, some MAR imputation methods can also be used for MCAR missing values and thus applied to proteomics data sets.

•Missing Not At Random (MNAR), which, on the contrary, has a targeted effect. In mass-spectrometry-based analysis, chemical species whose abundances are close enough to the limit of detection of the instrument record a higher rate of missing values. This is why MNAR-devoted imputation methods used in proteomics focus on left-censored data. (That is, the distribution of which with respect to the abundance is truncated on the left side, that is, on the region depicting the lower abundances.)

Second, the statistics literature contains numerous imputation methods devoted to MCAR or MAR, while very few are devoted to MNAR. The reason for this asymmetry is simple: Most of the MCAR/MAR mechanisms are generic to numerous application fields, so that it naturally focused statisticians' efforts. On the contrary, MNAR (including left-censored) mechanisms are discipline-specific, so that a precise understanding of the mechanism underlying the data generation is mandatory. This is why, in the comparisons depicted in ref 13, among the nine methods, only three MNAR-devoted approaches were considered, among which two are based on the same principle. Nonetheless, these nine methods have been compared on various data sets that are reported to have both MNAR and MCAR, yet in unknown proportions. As a result, even if a couple of MCAR/MAR devoted methods are shown to perform slightly better, it makes sense to wonder if this holds in general or if it is data-set-dependent.

Even though most of the conclusions of ref 13 are well supported, there is a need to consider the proportions of MCAR and MNAR as hidden variables. This idea is not new: Several recent works have proposed to perform imputation by estimating models (with maximum-likelihood<sup>14,15</sup> or with empirical Bayesian<sup>16</sup> methods), which are rich enough to account for both types of missingness mechanisms. To the best of our knowledge, no study has evaluated the behavior of an imputation method devoted to MNAR (respectively, devoted to MAR/MCAR) on a data set containing mainly MCAR (respectively, MNAR); however, this question is of prime importance to the practitioner, as it helps to guide the selection of an imputation algorithm according to the risk of corrupting the downstream analysis when using an unadapted imputation method.

In this work, we have considered real and simulated data sets on which MCAR and MNAR were introduced in controlled proportions and have compared the performances of various imputation methods. Numerous conclusions and recommendations can be drawn from these experiments; however, beyond them, our work pinpoints the fact that most of the conclusions regarding imputation methods cannot be claimed to hold in general. On the contrary, they should be contextualized according to each data set, the proportion of missing values, and their nature.

## 2. MATERIAL

### Simulated Quantitative Data Set

To generate artificial peptide abundance data, we used a simplified version of the model proposed in ref 5, which reads

$$y_{ij} = P_i + G_{ik} + \epsilon_{ij} \quad (1)$$

where  $y_{ij}$  is the log-transformed abundance of peptide  $i$  in the  $j$ th sample,  $P_i$  is the mean value of peptide  $i$ ,  $G_{ik}$  is the mean differences between the condition groups, and  $\epsilon_{ij}$  is the random error terms, which stands for the peptide-wise variance. Here  $P_i$  is randomly generated from a Gaussian distribution with mean  $\mu$  and standard deviation  $\sigma$ . The dynamic range of peptides (in logarithm scale) can be therefore approximated by  $[\mu - 3\sigma, \mu + 3\sigma]$ . We considered two groups  $k_1$  and  $k_2$  of replicates, for which  $P_i$  generation was conducted with  $\mu = 1.5$  and  $\sigma = 0.5$ . For each of the two groups, we selected two disjoint subsets of peptides (20% of the total number of peptides), and we added  $G_{ik}$  randomly drawn from the distribution previously mentioned to simulate a differential abundance between the peptides. Finally, the random error term has also been simulated by random draws from a Gaussian distribution with zero mean and standard deviation  $\sigma_e = 0.5$ . With these parameters, we simulated a log-transformed peptide abundance table with  $m = 1000$  peptides and  $n = 20$  replicates (equally split into groups  $k_1$  and  $k_2$ ).

To derive the protein abundance data, we have randomly generated a map describing the peptide/protein relationships by randomly drawing  $m$  integers from  $[1, m_{\text{prot}}]$ , where  $m$  is the number of peptides and  $m_{\text{prot}} < m$  is the number of proteins. ( $m_{\text{prot}}$  was set to  $m/2$ .)

### Real Quantitative Data Set

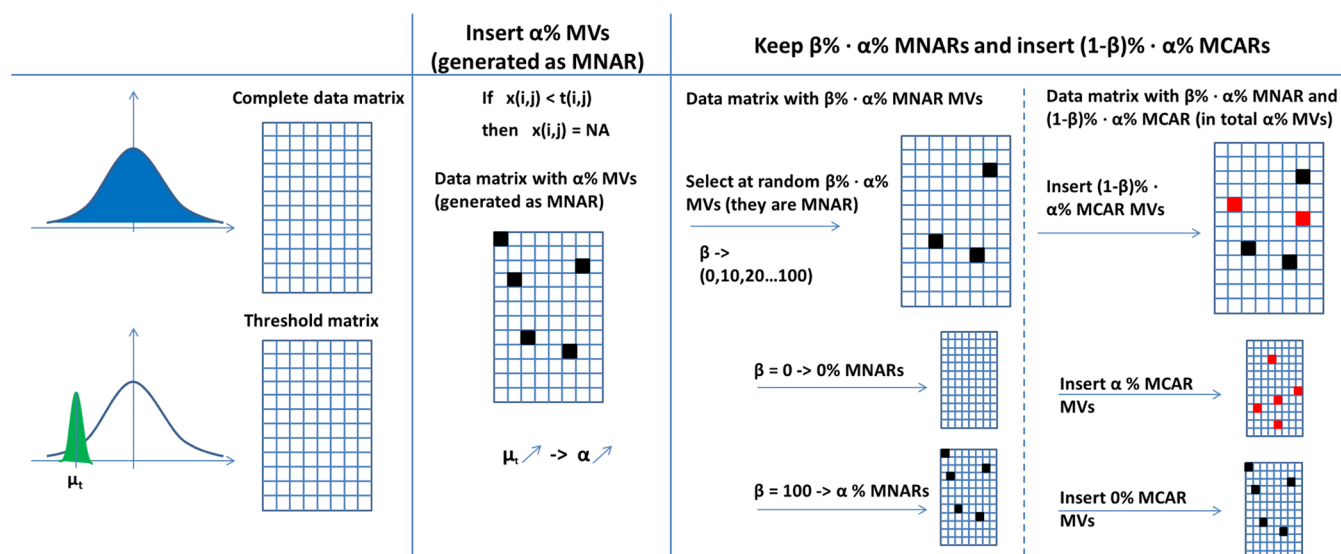
As a complement to the simulated data, we considered a real and publicly available data set, which has been collected during a study designed to compare human primary tumor-derived xenograph proteomes of the two major histological nonsmall cell lung cancer subtypes, adenocarcinoma (ADC) and squamous cell carcinoma (SCC), using Super-SILAC and label-free quantification.<sup>17</sup> The raw files were analyzed by MaxQuant (version 1.3.0.5). Peaks were searched against the UniProt human database (released July 2012; <http://www.uniprot.org>) using the Andromeda search engine included in MaxQuant. The data set within this package contains proteins intensity for six ADC and six SCC samples. The complete MaxQuant output file is available on the repository of the ProteomeXchange Consortium,<sup>18</sup> with the data set identifier PXD000438.

Because this study requires precisely controlling each missing values, one must work on a complete data set, that is, where no missing value shows up. This has been obtained from the raw peptide-level PXD000438 data set by filtering out the peptides that contain at least one missing value. Finally, the complete peptide-level matrix was log-transformed and median-normalized.

### MCAR and MNAR Incorporation

Let  $\alpha$  and  $\beta$  be the rate of missing values and the MNAR ratio, respectively. They read

$$\alpha = \frac{100 \cdot (\# \text{MNAR} + \# \text{MCAR})}{nm} \quad \beta = \frac{100 \cdot \# \text{MNAR}}{\# \text{MNAR} + \# \text{MCAR}} \quad (2)$$



**Figure 1.** Schematic view upon the strategy used for the missing data generation. This strategy allows us to control both for the total proportion of missing values generated as well as for the proportion of missing values, which are MNAR and MCAR.

For a given combination of  $\alpha$  and  $\beta$ , the missing values are incorporated in a complete data set as follows:

MNAR values are incorporated using a stochastic threshold, as follows: One randomly generates a threshold matrix  $T$  from a Gaussian distribution with parameters ( $\mu_t = q$ ,  $\sigma_t = 0.01$ ), where  $q$  is the  $\alpha$ th quantile of the abundance distribution in the complete quantitative data set. Then, each cell  $(i,j)$  of the complete quantitative data set is compared with  $T_{i,j}$ . If it is greater than or equal to  $T_{i,j}$ , the abundance is not censored. On the contrary, if it is strictly smaller than  $T_{i,j}$ , a Bernoulli draw with probability of success  $\frac{\beta \cdot \alpha}{100}$  determines if the abundance value is censored (success) or not (failure).

MCAR values are incorporated by replacing with a missing value the abundance value of  $nm \cdot \frac{(100 - \beta) \cdot \alpha}{100}$  randomly chosen cells in the table of the quantitative data set.

This strategy is summarized in Figure 1. We used it for any combination of values for  $\alpha \in [2\%, 52\%]$  and  $\beta \in [0\%, 100\%]$ .

### 3. METHODS

#### Imputation Algorithms

Because an exhaustive comparison of the missing value imputation algorithms is beyond the scope of this study, we selected a set of characteristic and widely applied methods, representing different families of imputation procedures and which are conceptually different. We considered:

- **kNN** ( $k$  Nearest Neighbors):<sup>3</sup> for a peptide showing missing values, the method consists of (i) finding  $k$  most similar peptides to the one considered (using a particular distance measure, e.g., Euclidean distance of Pearson's correlation coefficient) and (ii) imputing each missing value by averaging the  $k$  peptide values from the same replicate where that missing value occurred. Preliminary exploration of the range of parameter  $k$  showed that the imputation accuracy was rather stable for any  $k \in [10, 20]$  and reached its maximum at 11, so that we used this latter value.

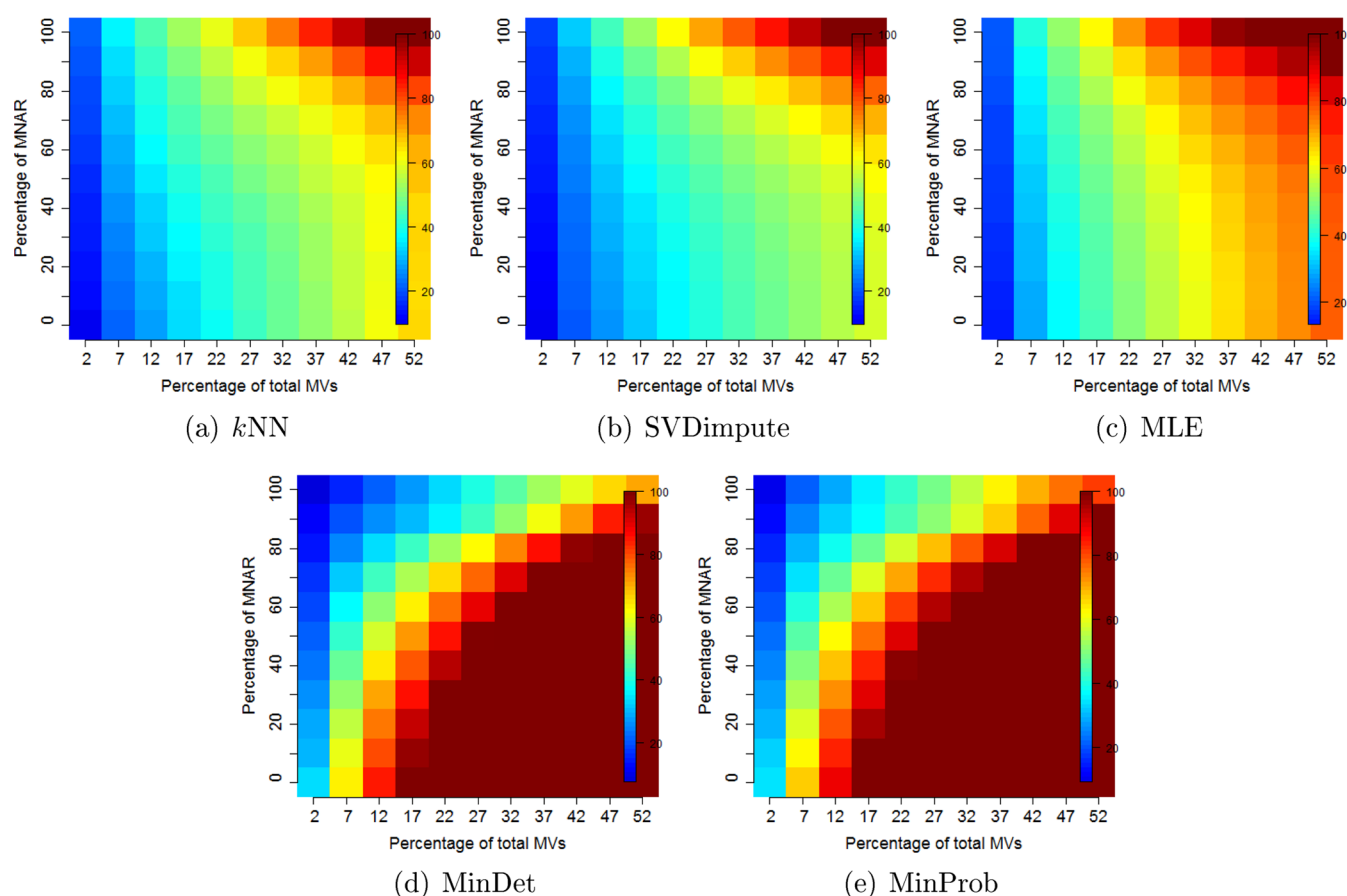
- **SVDimpute** (Imputation with Singular Value Decomposition):<sup>3</sup> The quantitative data set is considered a matrix on which mean centering and  $k$ -rank SVD are iteratively applied (where  $k \in [1, n/2]$ , where  $n/2$  is the number of replicates in a

given condition group), up to some convergence criterion. In our case,  $k$  was tuned to 1. ( $k = 1$  and 2 gave the greatest performances according to preliminary tests.)

- **MLE** (Imputation based on Maximum Likelihood Estimation): Assuming the quantitative data set obeys some law  $f_\theta$  of unknown parameter  $\theta$ , maximum likelihood estimation principle is used to derive an estimator  $\hat{\theta}$  of  $\theta$ , and missing values are then imputed by random draws of  $f_{\hat{\theta}}$ . The literature dedicated to missing value imputation based on MLE is vast, and we recommend refs 19 and 20 for a comprehensive survey of the topic. In this work, we employed the implementation available in the R package norm.<sup>21</sup>

- **MinDet** (Deterministic Minimum Imputation):<sup>22,23</sup> It simply replaces the missing values by the minimum value, either globally observed in the data set or observed in each sample. Here we used the  $10^{-4}$  quantile.

- **MinProb** (Probabilistic Minimum Imputation): It is a stochastic version of MinDet, so as to limit the bias introduced by multiple replacements with a unique value. The imputation is performed by replacing the missing values with random draws from a Gaussian distribution centered on the value used with MinDet and with a variance tuned to the median of the peptide-wise estimated variances.<sup>24</sup> We decided to focus on these five methods, as they represent well the various types of imputation methods: First, according to the taxonomies provided in refs 25 and 26, kNN, MinDet, and MinProb belong to the *prediction rules* methods, SVDimpute belongs to the *least-squares-based* methods, and finally, MLE belongs to the *maximum-likelihood-based* methods, so that this set of methods covers well the taxonomies of refs 25 and 26. Second, according to ref 13, MinDet and MinProb are *single value approaches*, kNN is a *local similarity approach*, and SVDimpute is a *global similarity approach*, so that the taxonomy of ref 13 is also covered. Third, MinDet and MinProb are designed to impute MNAR values, while kNN, SVDimpute, and MLE are designed for MCAR (and more generally MAR) values. Finally, MinDet is the most naive method to deal with MNAR (and often implemented as zero value imputation), while MLE and SVDimpute are particularly efficient on MCAR, so that comparing these three methods is insightful with regard to



**Figure 2.** RSR for the simulated quantitative data set; imputation is performed by considering: *k*NN (a), SVDimpute (b), MLE (c), MinDet (d), and MinProb (e).

the conclusions of ref 13 on the general dominance of MCAR/MAR-devoted methods. Let us also notice that no multiple imputation method is considered in our work, while, in practice, they provide the best results in the state-of-the-art. The reason is the following: Multiple imputation strategies amount to a boosting strategy, that is, the combination of several simple methods to stabilize the results; however, their behavior, efficiency, and adequation to the specificities of the data are directly related to those of the simple methods they are based on. As a result, we found it clearer to focus on the single imputation methods, so as to best describe and understand them, and to let the practitioner generalize our conclusion to multiple imputations. Finally, this set of algorithms has been chosen to represent a wide diversity of strategies, on which very general conclusions can be drawn.

#### Accuracy Measurements

In most of the experiments, the imputation step was followed by the aggregation of peptide abundances into protein abundances. (We estimated each protein abundance with the median abundance over the protein specific peptides.) However, in few specific experiments (see Section 4), the aggregation was conducted first (i.e., on peptide abundances that still contain missing values) and followed by imputation at the protein level.

In both cases, we evaluated the performances of the imputation algorithms in the same way: We considered the differences between the protein abundances in the original complete quantitative data set and in its counterpart containing missing values that have been imputed (either at protein or

peptide level). Such differences are classically summarized by the root-mean-square error (RMSE), yet many other variants exist.<sup>27</sup> Within our framework, we employed a normalized version of the RMSE called the RMSE-observations standard deviation ratio (RSR),<sup>28</sup> defined as follows

$$RSR(X_C, X_I) = \frac{RMSE(X_C, X_I)}{sd(X_C)} \quad (3)$$

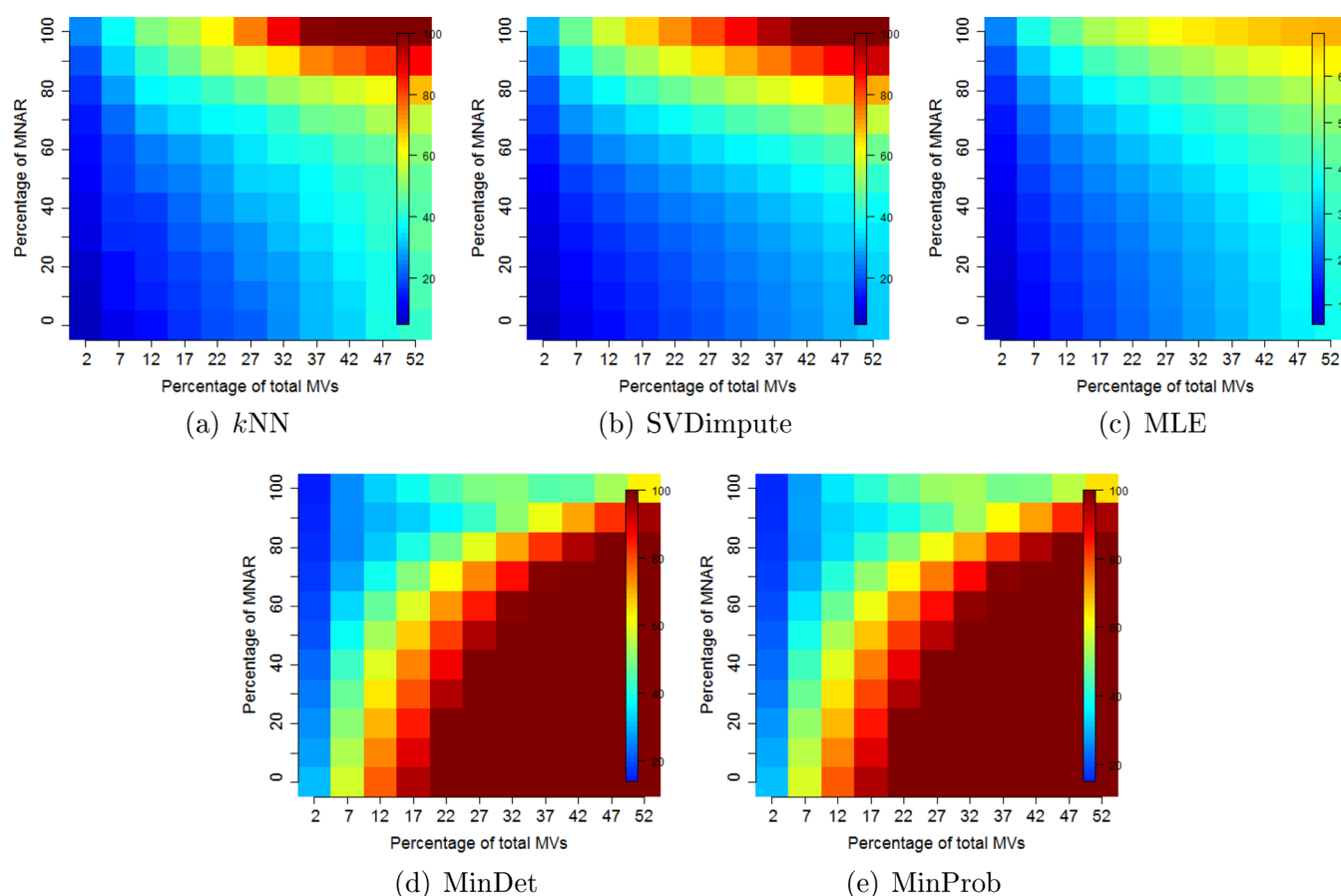
where  $X_C$  denotes the complete quantitative data set (before incorporating missing values), while  $X_I$  denotes the quantitative data set after the imputation of the missing values. The reported results corresponds to an average over 30 independent repetitions of the experiment (i.e., the random generation of missing values as well as their imputation, for a given tuning of  $\alpha$  and  $\beta$ ), so as to have more stable performance records.

## 4. RESULTS

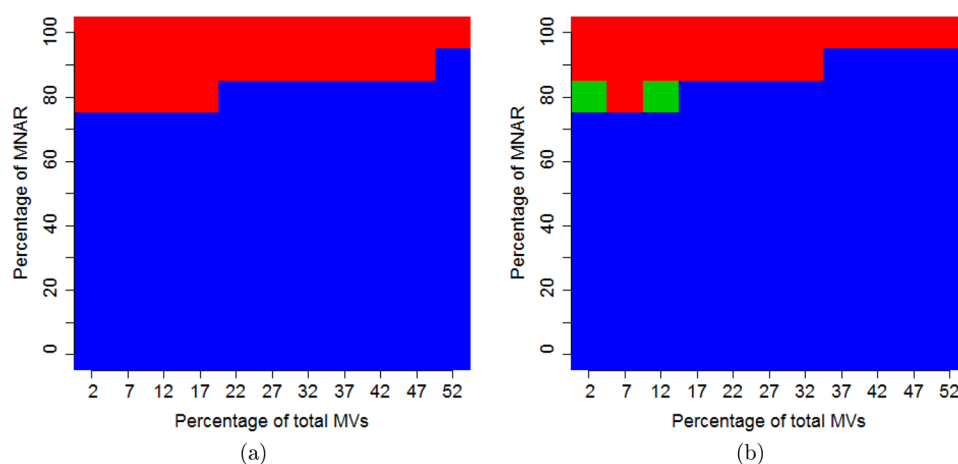
### MCAR-Devoted versus MNAR-Devoted Imputations

Figures 2 and 3 display a series of heatmaps (with a false color code, ranging from blue, which indicates low RSR, to red, which indicates high RSR) for the simulated and real data sets, respectively. Within each Figure, there are five graphics, corresponding to an imputation method each. Each heatmap displays the average performances (over 30 repetitions) of the imputation algorithm over the entire range of the experimental conditions (i.e., a proportion of missing values ranging from 2 to 52% and an MNAR ratio ranging from 0 to 100%). Several conclusions can be drawn from these Figures.





**Figure 3.** RSR for the real quantitative data set; imputation is performed by considering:  $k$ NN (a), SVDimpute (b), MLE (c), MinDet (d), and MinProb (e).



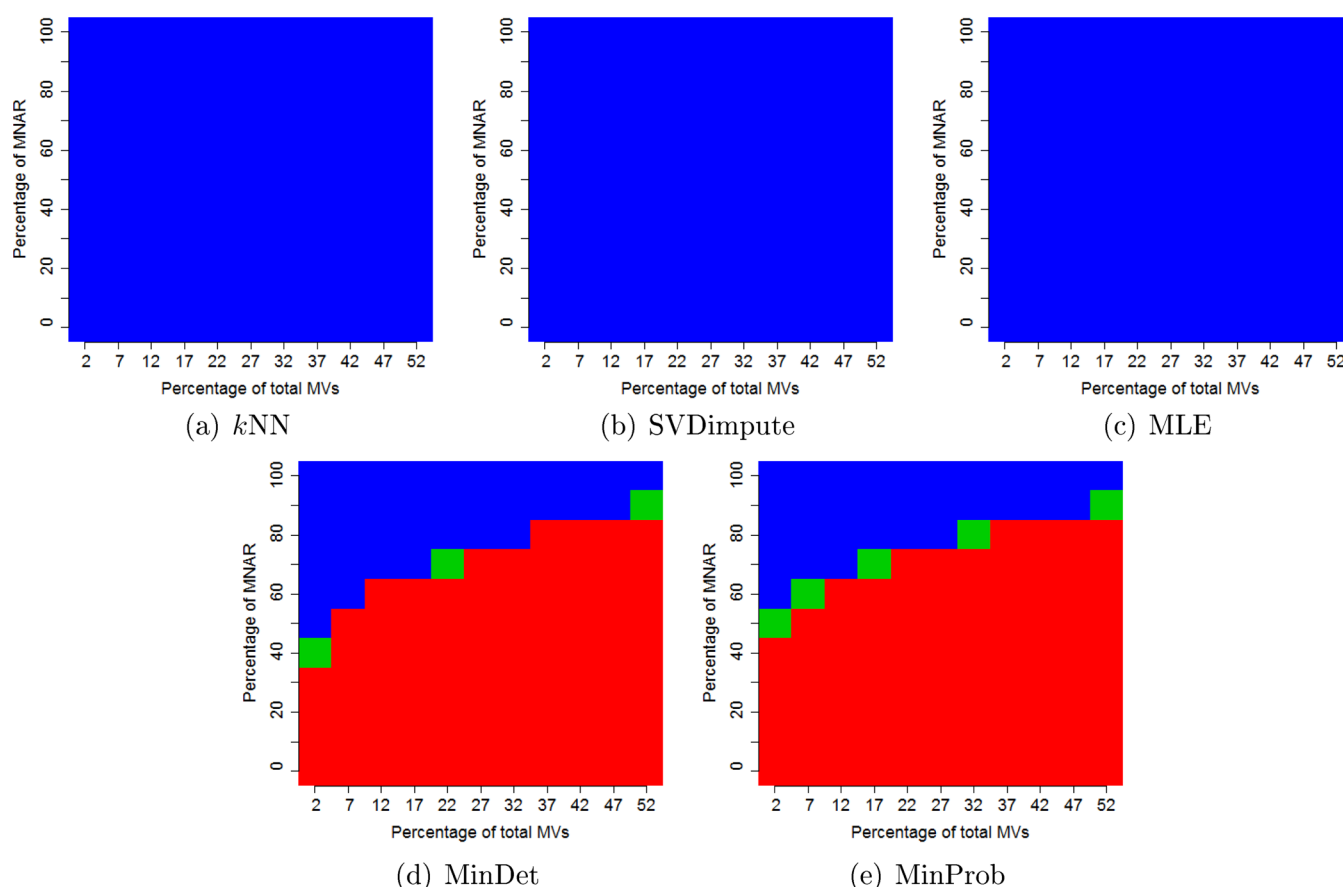
**Figure 4.** (a) Comparison of SVDimpute and MinDet on the simulated data set. (b) Comparison of MLE and MinDet on the real data set. Red color indicates an outperformance of MinDet, a blue color, an underperformance of MinDet, and a green color, a difference of performance that is not significant with a  $p$ -value threshold of 5%.

First, irrespective of the data set, all methods perform better when there are fewer missing values and become inaccurate with increasing proportion of missing values. Although expected, this result assesses the validity of our comparison protocol and of our simulations.

Second, two groups of algorithms can be identified with regard to the MNAR ratio: The first group is made of SVDimpute,  $k$ NN, and MLE, which perform better under a small MNAR ratio, while the second group, composed of

MinDet and MinProb, performs better under a larger MNAR ratio. This clearly indicates that depending on the nature of the majority of the missing values it is important to privilege either a MCAR/MAR-devoted method, such as advocated in ref 13, or, on the contrary, to favor a MNAR-devoted method, even if the latter is more naive and provide, on average, worse results.

Third, for each method, a similar behavior is observed on both the real and the simulated data sets. In the case of MinDet and MinProb, the similarity is almost perfect, with particular



**Figure 5.** Comparison of peptide-level and protein-level imputations for the simulated quantitative data set; imputation is performed by considering: *k*NN (a), SVDimpute (b), MLE (c), MinDet (d), and MinProb (e). Blue indicates peptide imputation superiority, red indicates protein imputation superiority, and green indicates a nonsignificant result (at 5% threshold).

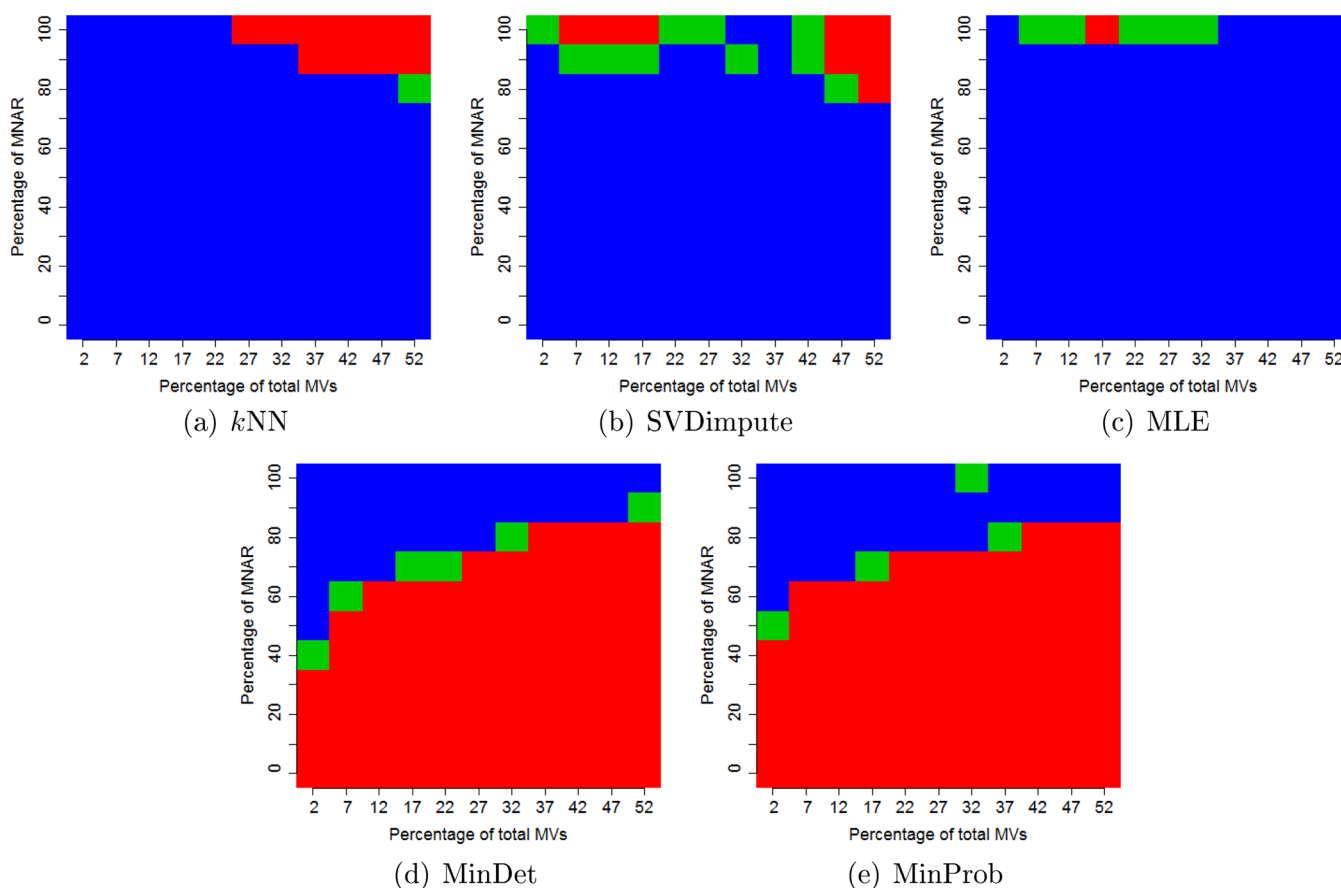
poor performance toward high percentages of nonrandom missing values (lower right corner). In the case of the three other methods, even if the similarity between the heatmaps derived from the real and simulated data sets is not as good, a pattern is well-conserved. In both cases, the best performance is reached with the lowest rate of missing value and the lowest MNAR ratio (lower left corner), while the worst performance is reached with the greatest rate of missing value and the greatest MNAR ratio (upper right corner). In addition, isoperformance lines are roughly parallel to an axis going from the upper left to the lower right corner. The global stability of this pattern indicates that, even if MCAR is possibly a simplistic process to account for the diverse nature of missing values that are not left-censored, the postulate at the root of these experiments is robust. Indeed, we postulated the strong influences of both (1) the rate of missing values and of the MNAR ratio as well as (2) the nature of the missing values to which a given imputation method is devoted.

Finally, if one averages the performances of the various imputation methods over all the experiments (which amounts to consider a mean color over each graphic), it appears that overall MCAR/MAR-devoted methods (SVDimpute, *k*NN and MLE) outperform MNAR methods (MinDet and MinProb). From this, we conclude that in the absence of any knowledge regarding the MNAR ratio (and assuming that all of the situations are equiprobable, which remains to be proven), it makes sense to privilege the former ones, such as advocated in ref 13.

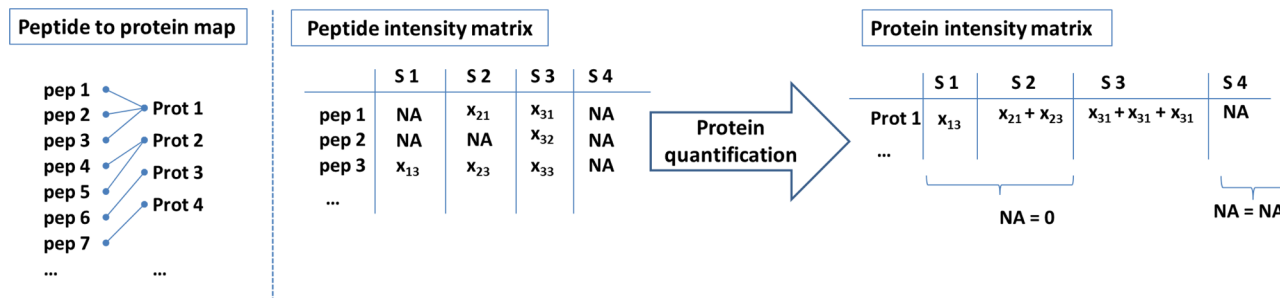
However, this averaging must not be overstated, as it is possible to show situations where even the worst MNAR method (MinDet) significantly outperforms the best MCAR/MAR methods (MLE or SVDimpute). To further demonstrate this, we applied an unpaired two-sample *t* test to assess the significance of the difference of accuracy, between the two following pairs of imputation methods: MinDet vs SVDimpute (for the simulated data set), and MinDet vs MLE (for the real data set). The results are reported in Figure 4. These comparisons demonstrate that when a high proportion (70% or more) of missing values are MNAR, MNAR imputation methods are preferred. Although such data sets are not widespread, they are not unheard of (see, for instance, refs 29 and 30), which advocates for the development of new methodologies that can estimate the nature of the majority of the missing values, so as to adapt the imputation method accordingly.

#### Peptide-Level versus Protein-Level Imputations

In the literature, there is no consensus on the preferred order with respect to missing values imputation and aggregation of peptide intensities into protein intensities. This is why both cases were considered in ref 13. We also repeated our experiments summarized in Figures 2 and 3, in a reversed context where the aggregation is performed first and the imputation is conducted at protein level. We have compared these two approaches using the methodology previously described and present the results of a significance analysis (at a *p*-value threshold of 5%) in Figures 5 and 6, where blue



**Figure 6.** Comparison of peptide-level and protein-level imputations for the real quantitative data set; imputation is performed by considering:  $k$ NN (a), SVDimpute (b), MLE (c), MinDet (d), and MinProb (e). Blue indicates peptide imputation superiority, red indicates protein imputation superiority, and green indicates a nonsignificant result (at 5% threshold).



**Figure 7.** Illustration of implicit missing value imputation during protein quantification from peptide intensity. Here the protein quantification is considered to be performed by summing the signal intensities of all peptides per protein.

indicates peptide imputation superiority, red indicates protein imputation superiority, and green indicates a nonsignificant result.

As illustrated by a high proportion of blue, peptide-level imputation is most of the time more accurate. Nevertheless, a major argument for protein-level imputation is the presence of fewer missing values; indeed, if several peptides are aggregated into a protein, this aggregation does not lead to a missing value, unless all of the peptide intensities are missing, so that numerous missing values are implicitly imputed by a value that is a neutral element with respect to the aggregation. For instance, in the case where

- protein intensities result from the sum of the peptide intensities (such as in ref 31); then, missing peptide

intensities do not contribute to the sum, so that the result is the same as if peptide missing values were imputed by zero

- protein intensities result from the mean of the peptide intensities (such as in ref 32); then, missing peptide intensities do not contribute to the mean, so that the result is the same as if peptide missing values were imputed by the mean value of the peptide intensities
- protein intensities result from a maximum function of the peptide intensities (sum or mean over the three most abundant peptides, maximum peptide abundance, and so on, such as in refs 33 and 34); then, the result is the same as if peptide missing values were imputed by zero or any other small intensity

For more general protein aggregation methods, based on more sophisticated functions (such as, for instance, weighted mean), the issue is the same (even if the formula of the neutral element may be less trivial). The above observations are schematically described in Figure 7, where protein-level imputation is equivalent to (i) applying an implicit imputation method on some peptide-level missing values that is neither controlled nor evaluated; (ii) performing the aggregation itself; and (iii) explicitly imputing the few remaining missing values. Because the total number of imputed missing values (whether implicit or explicit) is the same, it is preferable to consider an explicit and well-justified imputation for all missing values, which amounts to impute at peptide level and concurs with the results of Figures 5 and 6.

However, from Figures 5 and 6, it seemingly appears that when the data contain up to ~60% of MNAR values, and if an MNAR-devoted imputation method has been chosen a priori, it is more efficient to impute at the protein level. This observation highlights that on MCAR data an implicit and suboptimal imputation is more efficient than an MNAR imputation method. Deriving this result on the basis of the aforementioned observation (Figures 5 and 6) requires several steps:

(1) During the aggregation process, several MCAR peptides are combined with observed peptide intensities (there is very little chance that, assuming MCAR data, all peptides of a given protein are missing), leading to protein intensities rather than missing values.

(2) As opposed to (1), let us note that MNAR peptides correspond to genuine low abundance ions, so that there are good chances that one aggregates only missing values, leading to a missing value at the protein level.

(3) As a result from (1), it appears that if one has chosen to use an MNAR-devoted method, MCAR are either imputed by an unadapted method (at the peptide level), or implicitly imputed by the aggregation.

(4) As a result from (2), if one uses the same MNAR-devoted method, MNAR are roughly imputed in the same way, both at peptide and at protein levels.

(5) As a result from (3) and (4), one derives that the difference in the overall quality of the imputation (between peptide level and protein level imputation with an MNAR-devoted method) mainly relies on that of MCAR data.

(6) Let us now recall the original observation: "When the data contain up to ~60% of MNAR values, and if an MNAR-devoted imputation method has been chosen a priori, it is more efficient to impute at the protein level."

(7) Then, on the basis of (5) and (6), the observed difference in the overall comparison is mainly explained by the performances of the imputation on the 40% or less remaining MCAR values.

(8) From (6) and (7), one derives that on these MCAR values implicit protein-level imputation gives more accurate results.

(9) Then, combining (8) and (2) leads to the aforementioned conclusion: On MCAR, an implicit and suboptimal imputation is more efficient than a MNAR-devoted method.

As here, the implicit imputation of the aggregation is equivalent to a mean imputation (which can be seen as a poor MCAR method), it highlights that a bad MCAR method is more efficient on MCAR data than a good MNAR method. While this conclusion may appear trivial, it, however, stresses that the adequation between the nature of the missing values and the imputation strategy is more important than the

theoretical performances (i.e., regardless the nature of missing values) of the imputation algorithm.

In addition, a last conclusion can be drawn: Because the implicit imputation performed during the aggregation mainly operates on MCAR, so that mainly MNAR remains at the protein level, our results support the idea that the MNAR ratio is generally more important at the protein level than at the peptide level (such as observed in refs 29 and 30, for instance); however, this last conclusion must be cautiously interpreted. Indeed, it does not mean that if there is a lot of MNAR it is better to work at the protein-level: To derive such a conclusion, one would need to have mainly red cells in the upper lines of Figures 5 and 6 graphics; yet it holds only for a couple of them (Figures 6a,b), so that no general conclusion can be drawn.

Of course, if one changes the aggregation method, the comparison between peptide-level and protein-level imputations will lead to slightly different results, and we do not pretend to be exhaustive; however, even if the aggregation strategy is more elaborate than the three aforementioned ones (sum, mean, or max), the conclusions are of the same spirit: Whatever the aggregation function, it is most likely to have a neutral element that will act as the implicit imputation value, on the basis of which most of the aforementioned conclusions are elaborated.

## 5. CONCLUSIONS

Let us first summarize the conclusions of this work into four points. (1) Imputation should be performed at the peptide level because aggregating peptides into proteins beforehand amounts to performing a first implicit and, in most of the cases, suboptimal imputation. (2) In the absence of knowledge about the nature(s) of missing values in a particular quantitative proteomics data set, it makes sense to rely on a MCAR/MAR imputation method. This is supported by numerous experiments, including ours as well as those from ref 13 but also by theoretical arguments: By definition, missing values that should be imputed by small intensities can also show up in a MCAR context (so that they can also be imputed to some extent by MCAR-devoted imputation methods), while, on the contrary, a method devoted to left-censored missing value will systematically perform poorly on other types of missing values. (3) However, this conclusion should be moderated by the observation that the superiority of MAR/MCAR-devoted methods only holds on the average and should be contextualized, as cases arise where MNAR-devoted methods perform better than MCAR-devoted ones. Similarly, it appears that choosing a method adapted to the nature of the missing values is more important than choosing a method itself, regardless of the nature of missing values. As a consequence, before any imputation, the practitioner should identify the main or most likely nature among the missing values in his/her quantitative data set and impute accordingly. (4) Finally, while MNAR are best imputed by specific methods, other missing values are well accounted for by MAR/MCAR-devoted methods. Because it is accepted that many types of missing values coexist in most of the quantitative data sets (see, for instance, ref 13), hybrid strategies (based on both MNAR- and MAR/MCAR-devoted methods) should be considered in the future.

These elements shed a new light on the directions that methodological research should follow with regards to missing value imputation in quantitative proteomics. MNAR-devoted methods, which are less numerous and have been less



investigated in the general field of statistics, remain a subject of likely improvements. Concomitantly, important room is left to develop diagnosis tools, which are capable of categorizing the missing values according to the mechanism that generated them. This diagnosis can operate at different levels: (i) at the data set level, so that the imputation strategy is applied conditionally to the majority nature of missing values in the entire data set; (ii) at the peptide level, so that all of the missing values within a same peptide (in a given group of replicates) are assumed to be of a same nature; and (iii) at the missing value level, so as to have a most refined categorization of the missing values across the data set. Finally, once such diagnosis tools are available, it will be possible to elaborate hybrid strategies that process each group of missing values according to its nature, so as to best preserve the biological relevance of the quantitative data sets and of the biological conclusions.

## AUTHOR INFORMATION

### Corresponding Author

\*E-mail: [thomas.burger@cea.fr](mailto:thomas.burger@cea.fr).

### Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

This work was supported by the following funding: ANR-2010-GENOM-BTV-002-01 (Chloro-Types), ANR-10-INBS-08 (ProFI project, “Infrastructures Nationales en Biologie et Santé”, “Investissements d’Avenir”), EU FP7 program (Prime-XS project, Contract no. 262067), the Prospectom project (Mastodons 2012 CNRS challenge), and the BBSRC Strategic Longer and Larger grant (Award BB/L002817/1).

## REFERENCES

- (1) Stead, D. A.; Paton, N. W.; Missier, P.; Embury, S. M.; Hedeler, C.; Jin, B.; Brown, A. J. P.; Preece, A. Information quality in proteomics. *Briefings Bioinf.* **2007**, *9* (2), 174–188.
- (2) Albrecht, D.; Knemeyer, O.; Brakhage, A. A.; Guthke, R. Missing values in gel-based proteomics. *Proteomics* **2010**, *10* (6), 1202–1211.
- (3) Troyanskaya, O.; Cantor, M.; Sherlock, G.; Brown, P.; Hastie, T.; Tibshirani, R.; Botstein, D.; Altman, R. B. Missing value estimation methods for dna microarrays. *Bioinformatics* **2001**, *17* (6), 520–525.
- (4) Rubin, D. B. Inference and missing data. *Biometrika* **1976**, *63* (3), 581–592.
- (5) Karpievitch, Yuliya; Dabney, Alan; Smith, Richard Normalization and missing value imputation for label-free lc-ms analysis. *BMC Bioinf.* **2012**, *13* (Suppl. 16:S5), 1–9.
- (6) Hastie, T.; Tibshirani, R.; Narasimhan, B.; Chu, G. *Impute: Imputation for Microarray Data*. R package, version 1.42.0.
- (7) Lazar, C. *imputeLCMD: A Collection of Methods for Left-Censored Missing Data Imputation*. R package, version 2.0.
- (8) Stacklies, W.; Redestig, H.; Scholz, M.; Walther, D.; Selbig, J. *pcamethods - a bioconductor package providing pca methods for incomplete data*. *Bioinformatics* **2007**, *23* (9), 1164–1167.
- (9) Gatto, L.; Lilley, K. S. *Msnbase-an r/bioconductor package for isobaric tagged mass spectrometry data visualization, processing and quantitation*. *Bioinformatics* **2012**, *28* (2), 288–289.
- (10) Nie, L.; Wu, G.; Brockman, F. J.; Zhang, W. Integrated analysis of transcriptomic and proteomic data of *Desulfovibrio vulgaris*: zero-inflated poisson regression models to predict abundance of undetected proteins. *Bioinformatics* **2006**, *22* (13), 1641–1647.
- (11) Torres-García, W.; Zhang, W.; Runger, G. C.; Johnson, R. H.; Meldrum, D. e R Integrative analysis of transcriptomic and proteomic data of *Desulfovibrio vulgaris*: a non-linear model to predict abundance of undetected proteins. *Bioinformatics* **2009**, *25* (15), 1905–1914.
- (12) Torres-García, W.; Brown, S. D.; Johnson, R. H.; Zhang, W.; Runger, G. C.; Meldrum, D. R Integrative analysis of transcriptomic and proteomic data of *Shewanella oneidensis*: missing value imputation using temporal datasets. *Mol. BioSyst.* **2011**, *7* (4), 1093–1104.
- (13) Webb-Robertson, B.-J. M.; Wiberg, H. K.; Matzke, M. M.; Brown, J. N.; Wang, J.; McDermott, J. E.; Smith, R. D.; Rodland, K. D.; Metz, T. O.; Pounds, J. G.; Waters, K. M.; et al. Review, evaluation, and discussion of the challenges of missing value imputation for mass spectrometry-based label-free global proteomics. *J. Proteome Res.* **2015**, *14* (5), 1993–2001.
- (14) Karpievitch, Yuliya; Stanley, Jeff; Taverner, T.; Huang, J.; Adkins, J. N.; Ansong, C.; Heffron, F.; Metz, T. O.; Qian, W.-J.; Yoon, H.; Smith, R. D.; Dabney, A. R. A statistical framework for protein quantitation in bottom-up ms-based proteomics. *Bioinformatics* **2009**, *25* (16), 2028–2034.
- (15) Ryu, S. Y.; Qian, W.-J.; Camp, D. G.; Smith, R. D.; Tompkins, R. G.; Davis, R. W.; Xiao, W. Detecting differential protein expression in large-scale population proteomics. *Bioinformatics* **2014**, *30* (19), 2741–2746.
- (16) Koopmans, F.; Cornelisse, L. N.; Heskes, T.; Dijkstra, T. M. H. Empirical bayesian random censoring threshold model improves detection of differentially abundant proteins. *J. Proteome Res.* **2014**, *13* (9), 3871–3880.
- (17) Zhang, W.; Wei, Y.; Ignatchenko, V.; Li, L.; Sakashita, S.; Pham, N.-A.; Taylor, P.; Tsao, M. S.; Kislinger, T.; Moran, M. F. Proteomic profiles of human lung adeno and squamous cell carcinoma using super-silac and label-free quantification approaches. *Proteomics* **2014**, *14* (6), 795–803.
- (18) Vizcaino, J. A.; Deutsch, E. W.; Wang, R.; Csordas, A.; Reisinger, F.; Ríos, D.; Dianes, J. A.; Sun, Z.; Farrah, T.; Bandeira, N.; et al. Proteomexchange provides globally coordinated proteomics data submission and dissemination. *Nat. Biotechnol.* **2014**, *32* (3), 223–226.
- (19) Ibrahim, J. G.; Chen, M.-H.; Lipsitz, S. R.; Herring, A. H. Missing-data methods for generalized linear models. *J. Am. Stat. Assoc.* **2005**, *100* (469), 332–346.
- (20) Schafer, J. L.; Graham, J. W. Missing data: Our view of the state of the art. *Psychological Methods* **2002**, *7* (2), 147–177.
- (21) Schafer, J. L. *NORM: Analysis of Incomplete Multivariate Data under a Normal Model*, 3rd ed.; The Methodology Center, The Pennsylvania State University: University Park, PA, 2008.
- (22) Almeida, J. S.; Stanislaus, R.; Krug, E.; Arthur, J. M. Normalization and analysis of residual variation in two-dimensional gel electrophoresis for quantitative differential proteomics. *Proteomics* **2005**, *5* (5), 1242–1249.
- (23) Meleth, S.; Deshane, J.; Kim, H. The case for well-conducted experiments to validate statistical protocols for 2d gels: different pre-processing = different lists of significant proteins. *BMC Biotechnol.* **2005**, *5* (1), 7.
- (24) Chich, J.-F.; David, O.; Villers, F.; Schaeffer, B.; Lutowski, D.; Huet, S. Statistics for proteomics: Experimental design and 2-de differential analysis. *J. Chromatogr. B: Anal. Technol. Biomed. Life Sci.* **2007**, *849* (1–2), 261–272.
- (25) Wasito, I.; Mirkin, B. Nearest neighbour approach in the least-squares imputation algorithms. *Inf. Sci.* **2005**, *169*, 1–25.
- (26) Little, R. J. A. Regression with missing x’s: A review. *J. Am. Stat. Assoc.* **1992**, *87* (420), 1227–1237.
- (27) Oh, S.; Kang, D. D.; Brock, G. N.; Tseng, G. C. Biological impact of missing-value imputation on downstream analyses of gene expression profiles. *Bioinformatics* **2011**, *27* (1), 78–86.
- (28) Chen, H.; Xu, C.-Y.; Guo, S. Comparison and evaluation of multiple gcms, statistical downscaling and hydrological models in the study of climate change impacts on runoff. *J. Hydrol.* **2012**, *434*–435 (0), 36–45.
- (29) Ferro, M.; Brugière, S.; Salvi, D.; Seigneurin-Berny, D.; Court, M.; Moyet, L.; Ramus, C.; Miras, S.; Mellal, M.; Le Gall, S.; Kieffer-Jaquinod, S.; et al. At\_chloro, a comprehensive chloroplast proteome database with subplastidial localization and curated information on envelope proteins. *Mol. Cell. Proteomics* **2010**, *9* (6), 1063–1084.

(30) Tomizioli, M.; Lazar, C.; Brugière, S.; Burger, T.; Salvi, D.; Gatto, L.; Moyet, L.; Breckels, L. M.; Hesse, A.-M.; Lilley, K. S.; et al. Deciphering thylakoid sub-compartments using a mass spectrometry-based approach. *Mol. Cell. Proteomics* **2014**, *13* (8), 2147–2167.

(31) Roulhac, P. L.; Ward, J. M.; Thompson, J. W.; Soderblom, E. J.; Silva, M.; Moseley, M. A.; Jarvis, E. D. Microproteomics: Quantitative proteomic profiling of small numbers of laser-captured cells. *Cold Spring Harbor Protocols* **2011**, *2011* (2), 218–234.

(32) Ludwig, C.; Claassen, M.; Schmidt, A.; Aebersold, R. Estimation of absolute protein quantities of unlabeled samples by selected reaction monitoring mass spectrometry. *Mol. Cell. Proteomics* **2012**, *11* (3), M111.013987.

(33) Grossmann, J.; Roschitzki, B.; Panse, C.; Fortes, C.; Barkow-Oesterreicher, S.; Rutishauser, D.; Schlapbach, R. Implementation and evaluation of relative and absolute quantification in shotgun proteomics with label-free methods. *J. Proteomics* **2010**, *73* (9), 1740–1746.

(34) Silva, J. C.; Gorenstein, M. V.; Li, G.-Z.; Vissers, J. P. C.; Geromanos, S. J. Absolute quantification of proteins by lcms: A virtue of parallel ms acquisition. *Mol. Cell. Proteomics* **2005**, *5* (1), 144–156.