
Principal Components Analysis based frameworks for efficient missing data imputation algorithms

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Missing data is a commonly occurring problem in practice. Many imputation meth-
2 ods have been developed to fill in the missing entries. However, not all of them
3 can scale to high-dimensional data, especially the multiple imputation techniques.
4 Meanwhile, the data nowadays tending toward high-dimensional. Therefore, in this
5 work, we propose *Principal Component Analysis Imputation* (PCAI), a simple but
6 versatile framework based on Principal Component Analysis (PCA) to speed up
7 the imputation process and alleviate memory issue of many available imputation
8 techniques, without sacrificing the imputation quality in term of MSE. In addition,
9 the frameworks can be used even when some or all of the missing features are
10 categorical, or when the number of missing features is large. Next, we introduce
11 *PCA Imputation - Classification* (PIC), an application of PCAI for classification
12 problem with some adjustment. We validate our approach by experiments on vari-
13 ous scenarios, which shows that PCAI and PIC can work with various imputation
14 algorithms, including the state-of-the-art ones and improve the imputation speed
15 significantly, while achieving competitive mean square error/classification accuracy
16 compared to direct imputation (i.e., impute directly on the missing data).

17 1 Introduction

18 Despite recent efforts in directly handling missing data [1, 2, 3, 4], missing data imputation approaches
19 [5, 6, 7] remain commonly used. This is because directly handling missing data can be complicated
20 and usually are developed for specific target problems or models, while imputation can be more
21 versatile. Specifically, an important advantage of imputation is that the imputed data becomes
22 *complete*, i.e., no longer have any missing values. Therefore, it is easier to continue with other
23 preprocessing steps, analysis, and data visualizations. Furthermore, one can deploy many models
24 and choose the best one by using the available packages or software tools for non-missing data.
25 Meanwhile, directly handling missing data strategies do not have these advantages. They are more
26 complicated and not that readily available.

27 Many techniques have been developed for missing data imputation, ranging from traditional tech-
28 niques such as MICE [5], K-Nearest Neighbors to recent machine learning/deep learning techniques
29 such as GAIN [6], DL-GSA [7]. However, most of them are computationally expensive for big
30 datasets. For example, experiments in [8] show that under their experiment settings, for Fashion
31 MNIST [9], a dataset of 70,000 samples and 784 features, the MICE [5] and missForest [10] tech-
32 niques are unable to finish the imputation process within three hours for a missing rate (the ratio
33 between the number of missing entries versus the total number entries in the dataset) of 20%. Since
34 datasets nowadays are trending towards larger sizes [11], with hundreds of thousands of features
35 [12], it is crucial to speed up the available imputation techniques. Taking into account resource

36 consumption and availability such speed up cannot be achieved by only providing more and better
37 hardware but by the development of new methods.

38 To achieve this goal, this work introduces two novel frameworks based on Principal Component
39 Analysis (PCA) to speed up the imputation process of many available techniques or the imputation-
40 classification process for missing data classification problems. The first framework, **PCA Imputation**
41 **(PCAI)** is proposed to speed up the imputation speed by partitioning the data into the fully observed
42 features partition and the partition of features with missing data. After that, the imputation of the
43 missing part is performed based on the union of the PCA - reduced version of the fully observed
44 part and the missing part. Interestingly, it turns out that the method has a great potential to aid the
45 performance of methods that rely on many parameters, such as Deep Learning imputation techniques.
46 Meanwhile, the second one, **PCA Imputation - Classification (PIC)** is proposed to deal with the
47 missing data classification problems where dimension reduction is desirable in advance of the model
48 training step. PIC is based on PCAI with some modifications. Note that these frameworks are
49 different from the methods developed for principal component analysis under missing data presented
50 in [13, 14], which are about how to conduct PCA when the data contains missing values.

51 In summary, the contributions of this article are: (i) we introduce **PCAI** to improve the imputation
52 speed of many available imputation techniques; (ii) we introduce **PIC** to deal with missing data
53 classification problems where dimension reduction is desirable; (iii) we analyze the potential strength
54 and drawbacks of these approaches; and (iv) we illustrate via experiments that our frameworks
55 can work with various imputation strategies while achieve comparable or even lower mean square
56 error/higher classification accuracies compared to the corresponding original approaches, and alleviate
57 the memory issue in some approaches.

58 The rest of the paper is organized as follows. In Section 2 and Section 3, we review some related
59 work in the field of missing data, and review two popular formulations of PCA. Next, in Section
60 4, Section 5, and Section 6, we introduce our novel PCAI and PIC frameworks, and study their
61 relation to previous works, respectively. After that, in Section 7, we demonstrate their capabilities
62 via experiments on various datasets. The paper ends with conclusions, remarks, and future works in
63 Section 8.

64 2 Related Works

65 Various works have been published on missing data imputation to deal with different data analysis
66 situations. As an example, if one is interested in modeling the uncertainty associated with the
67 imputation, suitable approaches can be multiple or Bayesian imputation techniques such as multiple
68 imputations using Deep Denoising Autoencoders [15], Bayesian Principal Component Analysis-
69 based imputation [16], and extreme learning machine multiple imputation [17]. In addition, graphical
70 models can be prominent candidates when transparency, estimability, and testability are desirable, and
71 these approaches can provide meaningful performance guarantees even if the missing values are not at
72 random [18]. Next, for continuous data, matrix completion techniques such as Fast Alternating Least
73 Squares [19], softImpute [20] can quickly give good results. In biology, the missing values are often
74 categorical, and the imputed values need to be interpretable. In such cases, classification techniques
75 or tree-based methods such as decision trees and fuzzy clustering with iterative learning (DIFC) [21],
76 missForest [10], the DMI algorithm [22], and sequential regression trees [23] are well-suited. In
77 addition, some recently developed methods that can handle mixed missing data are SICE [24], FEMI
78 [25], and HCMM-LD [26]. When the sample sizes are large enough compared to the number of
79 features, deep learning techniques such as Multiple Imputation Using Deep Denoising Autoencoders
80 [15], DL-GSA [7], and Swarm Intelligence-Deep Neural Network [27] can be powerful imputers.
81 However, it is worth noting that deep learning methods usually require more data than statistical
82 imputation approaches. Some other popularly used missing data imputation methods are multiple
83 imputation by chained equation (MICE) [5], K-nearest Neighbors imputation (KNNI) [28], and mean
84 imputation [28].

85 In addition, for the purpose of data imputation and data type, for classification, the impact of
86 imputation techniques on different classifiers may vary. Specifically, [28] compares the performance
87 of logistic regression with regularization, k-nearest neighbours (kNN), random forest, classification
88 tree, and xgboost classifiers [28] on datasets with missing entries. They use different imputation
89 methods (mean imputation/ MICE imputation [5]/ missForest [10]/ random imputation/ softImpute

90 [20]/ hot deck imputation, kNN imputation) and compare the performance. According to the paper,
 91 mean imputation seems to outperform other counterparts for logistic regression with regularization
 92 and kNN, random imputation wins for random forest, missForest seems to be the best imputer for
 93 classification tree, and hot deck imputation is the best for xgboost.

94 With the rapid growth of data size [11, 12], it is necessary to speed up the available imputation
 95 methods because many current approaches remain too slow for big datasets, as pointed out in an
 96 example in Section 1. This is where a popular dimension reduction method like PCA can come to use.
 97 PCA projects the original higher-dimensional dataset into a representation of lower dimensionality
 98 by extracting and retaining important information from the data and expressing this new information
 99 based on a set of orthogonal vectors known as principal components. Its goal is to find linear
 100 transformations of the original data that retain the maximal amount of variance. Note that there are
 101 some works on PCA under data missingness. For example, [13] considers the problem of finding
 102 principal components as an optimization problem of an objective function and proposes iterative
 103 solutions to it. On the other hand, [29] proposes a multiple imputation method for the estimates of the
 104 parameters (components and axes) of PCA to take into account the variability due to missing values.
 105 However, our work is different from these works in the sense that they target the problem of how to
 106 perform PCA for a dataset with missing data. Meanwhile, our frameworks utilize PCA to speed up
 107 the imputation processor to reduce the ratio between the number of features and the sample size.

108 3 Preliminaries

109 Let $\mathbf{X} = [x_{ij}]$ where $i = 1, \dots, n; j = 1, \dots, p$ be a input data matrix of n samples, p features. In
 110 addition, assume that the features are centered and scaled. We review two popular formulations of
 111 PCA, which we refer to as PCA formulation 1 (**PCA-form1**) and PCA formulation 2 (**PCA-form2**).

112 3.1 PCA based on covariance matrix (PCA-form1)

113 Let Σ be the covariance matrix of \mathbf{X} . Next, let $(\lambda_1, \mathbf{v}_1), \dots, (\lambda_p, \mathbf{v}_p)$ be the sorted eigenvalue-
 114 eigenvector pairs of Σ such that $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$. Suppose that we choose the first r pairs
 115 for dimension reduction. Then the amount of variance explained by these r pairs is

$$\frac{\lambda_1 + \lambda_2 + \dots + \lambda_r}{\lambda_1 + \lambda_2 + \dots + \lambda_p} \quad (1)$$

116 In addition, let $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_r]$. Then the dimension reduced version of \mathbf{X} is \mathbf{XV} .

117 3.2 PCA based on the input matrix \mathbf{X} (PCA-form2)

118 The solution of PCA can also be produced based on the singular value decomposition of \mathbf{X} [30]:

$$\mathbf{X} = \mathbf{UDW}^T \quad (2)$$

119 where \mathbf{U} is an $n \times p$ orthogonal matrix, \mathbf{W} is a $p \times p$ orthogonal matrix, and \mathbf{D} is a $p \times p$ diagonal
 120 matrix whose diagonal elements are $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$. Suppose that r eigenvalues are used,
 121 then the projection matrix is $\mathbf{V} = \mathbf{W}_r \mathbf{W}_r^T$ where \mathbf{W}_r consists of the first r columns of \mathbf{W} . Then
 122 the dimension reduced version of \mathbf{X} is also \mathbf{XV} .

123 4 PCA Imputation (PCAI)

124 In this section, we detail our PCAI framework, a PCA based framework that is capable of significantly
 125 improving the imputation speed of an imputer for high dimensional data, alleviating the memory
 126 issue for many approaches.

127 To start with some notations, let $pca(A)$ be a function of a data matrix A . The function returns
 128 (\mathcal{R}_A, V) where \mathcal{R}_A is the PCA-reduced version of A , and V is the projection matrix where the i^{th}
 129 column of V is the eigenvector corresponding to the i^{th} largest eigenvalue. In addition, denote by
 130 $\mathcal{A} \cup \mathcal{B}$ the columnwise concatenation of two data partition \mathcal{A} and \mathcal{B} of relevant sizes. Next, suppose
 131 that we have a dataset $\mathcal{D} = \mathcal{F} \cup \mathcal{M}$, where \mathcal{F} consists of data from fully observed features and \mathcal{M}
 132 consists of data from features with missing values.

133 The framework is as depicted in Algorithm 1. We first conduct dimension reduction on the fully
 134 observed partition \mathcal{F} , which produces a reduced version \mathcal{R} of \mathcal{F} . Then, the imputation of \mathcal{M} is done
 135 on the set $\mathcal{R} \cup \mathcal{M}$ instead of $\mathcal{D} = \mathcal{F} \cup \mathcal{M}$ as how imputations are usually done (i.e., impute directly
 136 on the original missing data). In conducting dimension reduction, we expect to reduce the dimension
 137 of the fully observed partition so that the imputation of \mathcal{M} can be faster.

Algorithm 1 PCAI framework

Require:

- $\mathcal{D} = \mathcal{F} \cup \mathcal{M}$ where \mathcal{F} is the fully observed partition and \mathcal{M} is the partition with missing values
- Imputer I
- PCA algorithm pca

Procedure:

$(\mathcal{R}, V) \leftarrow pca(\mathcal{F})$
 $\mathcal{M}' \leftarrow$ the imputed version of \mathcal{M} based on $\mathcal{R} \cup \mathcal{M}$
return Imputed version \mathcal{M}' of \mathcal{M}

138 For the choice of the PCA formulation, note that if the number of samples is larger than the number
 139 of features in \mathcal{F} , then the size of the covariance matrix is smaller than the size of \mathcal{F} . Therefore, one
 140 may expect using the formulation of PCA based on the covariance matrix, as in Section 3.1, to be
 141 faster. Meanwhile, if the number of features in \mathcal{F} is larger than the sample size, then the covariance
 142 matrix of \mathcal{F} is larger than \mathcal{F} . Therefore, in such a case, it is better to use the PCA formulation based
 143 on the data itself, i.e., formulation as in Section 3.2.

144 One may reckon that using $\mathcal{R} \cup \mathcal{M}$ instead of $\mathcal{F} \cup \mathcal{M}$ may lead to loss of information due to
 145 dimension reduction and therefore lower the quality of imputation. However, as will be illustrated
 146 in the experiments, the differences between the mean squared error of the imputed version versus
 147 the ground truth for these approaches are only slightly different, and many times, PCAI seems to be
 148 slightly better. This is possibly because PCA retains the important information from the data while
 149 removing some noise, and therefore helps improving the imputation quality. However, PCAI also has
 150 some shortcomings. For problems where the sample size n is smaller than the number of features in
 151 the fully observed block q , if PCA-form1 is used, the covariance matrix has the size of $q \times q$, which
 152 is bigger than the size $n \times q$ of the fully observed partition \mathcal{F} . This may make the PCA dimension
 153 reduction process become computationally expensive, rendering PCAI to be slower than imputing
 154 directly on the original missing data. This issue will be illustrated in the experiment section.

155 5 PCAI for classification (PIC)

156 In this section, we discuss a straightforward application of PCAI in classification, with a slight
 157 modification for classification problems where it is desirable to conduct a dimension reduction before
 158 training a model, such as when the number of features is much larger than the sample size.

159 Since PCAI conducts PCA on the fully observed partition \mathcal{F} , it reduces the dimensions for a portion
 160 of the data. Therefore, rather than imputing values using the PCAI framework and then conducting
 161 a dimension reduction step on $\mathcal{F} \cup \mathcal{M}'$, one can perform dimension reduction on \mathcal{M}' to get \mathcal{R}' , a
 162 PCA-reduced version of \mathcal{M}' . Then, one can use $\mathcal{F} \cup \mathcal{R}'$ as reduced dimension data. As will be shown
 163 in the experiments, this speeds up the imputation and classification process significantly. This is the
 164 basic idea of our *Principle component Imputation for Classification (PIC)* framework.

165 PIC operates as shown in Algorithm 2. The procedure starts by performing PCA on the training fully
 166 observed partition \mathcal{F}_{train} , which gives the reduced version \mathcal{R}_{train} of \mathcal{F}_{train} and a projection matrix
 167 V . Next, we project \mathcal{F}_{test} on V to get the reduced version \mathcal{R}_{test} of \mathcal{F}_{test} . Then, we impute \mathcal{M}_{train}
 168 on $\mathcal{R}_{train} \cup \mathcal{M}_{train}$ to get the imputed version \mathcal{M}'_{train} . Next, we impute \mathcal{M}_{test} on $\mathcal{R}_{test} \cup \mathcal{M}_{test}$
 169 to get the imputed version \mathcal{M}'_{test} . After that, if $reduce_{miss}$ is set to true, we perform dimension
 170 reduction on $\mathcal{M}'_{train}, \mathcal{M}'_{test}$. Then, we train the classifier on $\mathcal{R}_{train} \cup \mathcal{R}'_{train}$, i.e., the union of the
 171 reduced version of \mathcal{F}_{train} and the reduced version of \mathcal{M}_{train} . For prediction of a vector $\mathbf{x} \in \mathcal{D}$,
 172 we can decompose \mathbf{x} into $\mathbf{x} = (\mathbf{x}_{\mathcal{F}}, \mathbf{x}_{\mathcal{M}})$. After that, we can project $\mathbf{x}_{\mathcal{F}}$ on V to get a projection \mathbf{r} .
 173 Similarly, we can project $\mathbf{x}_{\mathcal{M}}$ on V to get projection \mathbf{r}' . Finally, we can predict the label of \mathbf{x} using
 174 the classifier C with input $(\mathbf{r}, \mathbf{r}')$.

175 Note that $reduce_{miss}$ is an option. When the number of features in the missing partition \mathcal{M} is
 176 large, one may be interested in reducing the dimension of \mathcal{M}' , and therefore, set $reduce_{miss}$ to *True*.
 177 However, when the number of features in the missing partition is small, one may want to keep it to
 178 *False*. Also, since PIC is a straightforward application of PCAI for classification, the choice of PCA
 179 formulation should be used is similar to PCAI, which is analyzed in the previous section.

Algorithm 2 PIC framework

Require:

- $\mathcal{D} = \mathcal{F} \cup \mathcal{M}$ where \mathcal{F} is the fully observed partition and \mathcal{M} is the partition with missing values
- $reduce_{miss} = True/False$: if *True*, perform dimension reduction on the imputed partitions; if *False*, do not perform dimension reduction on the imputed partitions
- $\mathcal{F}_{train}, \mathcal{F}_{test}$: the training and testing data of the fully observed partition \mathcal{F} , respectively
- $\mathcal{M}_{train}, \mathcal{M}_{test}$: the training and testing data of the partition that has missing data \mathcal{M} , respectively

- Imputer I , classifier C , PCA algorithm pca

Procedure:

```

( $\mathcal{R}_{train}, V$ )  $\leftarrow$   $pca(\mathcal{F}_{train})$ 
 $\mathcal{R}_{test} \leftarrow \mathcal{F}_{test}V$ 
 $\mathcal{M}'_{train} \leftarrow$  imputed version of  $\mathcal{M}_{train}$  based on  $\mathcal{R}_{train} \cup \mathcal{M}_{train}$ 
 $\mathcal{M}'_{test} \leftarrow$  imputed version of  $\mathcal{M}_{test}$  based on  $\mathcal{R}_{test} \cup \mathcal{M}_{test}$ 
if  $reduce_{miss}$  then
  ( $\mathcal{R}'_{train}, W$ )  $\leftarrow$   $pca(\mathcal{M}'_{train})$ 
   $\mathcal{R}'_{test} \leftarrow \mathcal{M}'_{test}W$ 
  Train the classifier  $C$  based on  $\mathcal{R}_{train} \cup \mathcal{R}'_{train}$ 
  Classify based on  $\mathcal{R}_{test} \cup \mathcal{R}'_{test}$ ,
else
  Train the classifier based on  $\mathcal{R}_{train} \cup \mathcal{M}'_{train}$ 
  Classify based on  $\mathcal{R}_{test} \cup \mathcal{M}'_{test}$ 
end if
return trained classifier  $C$ 

```

180 **6 Relation to previous works**

181 Various works have been done on PCA that are related to missing data, which mostly can be
 182 categorized into missing values imputation using PCA, or dimension reduction using PCA under
 183 missing values. Some typical works that make use of PCA for missing values imputation are
 184 probabilistic PCA for missing flow volume data imputation [31]; chunk-wise iterative PCA for
 185 data imputation on datasets with many samples[32]; [14] proposes a fast algorithm for PCA under
 186 missing data that help in case of sparse, high dimensional data; [33] analyze maximum likelihood
 187 PCA (MLPCA) on maximum likelihood missing data imputation; and [34] proposed an imputation
 188 approach based on PCA and factorial analysis for mixed data.

189 Next, PCA under missing values was first studied in [35], where only one component and one
 190 imputation iteration are used. After that, [36] proposes a method based on MLPCA, where the
 191 method assigns large variance to missing values prior to implementing the method, which aim to
 192 guide the algorithm to fit a PCA model disregarding those points. Also, [37] introduce EM algorithm
 193 for building a PCA model that can deal with missing data. More recently, [38] proposes new
 194 techniques for building a PCA model with missing data: known data regression (KDR), projection to
 195 the model plane, KDR with principal component regression. In addition, [39] studies estimation and
 196 imputation in Probabilistic PCA when the data is missing not at random.

197 Different from the previous approaches, PCAI is a framework to speed up the imputation process,
 198 which can be used with various imputation methods, including the aforementioned PCA imputation
 199 algorithms and the state-of-the-art imputation algorithms such as softImpute [20], MissForest [10],
 200 GAIN [6]. In addition, note that since PCAI and PIC conduct dimension reduction on the fully
 201 observed partition \mathcal{F} , and not the missing portion \mathcal{M} if $reduce_{miss} = False$, they can handle missing
 202 data even if categorical features presents in the missing portion \mathcal{M} , when being used with imputers
 203 that's capable of handling categorical/mixed data (MissForest [10], SICE [24], FEMI [25], etc.). In

Table 1: Description of datasets used in our experiments

Dataset	# Classes	# Features	# Samples
Parkinson [42]	2	754	756
Fashion MNIST [9]	10	784	70000
Gene [43]	5	20531	801

204 addition, even if there exists categorical and continuous features in \mathcal{M} ; or $reduce_{miss} = True$ and
 205 there exists categorical and continuous features in \mathcal{M} , one can easily adjust the algorithm to conduct
 206 PCA on continuous features only. The previously mentioned PCA based approaches are, however,
 207 can only be used for continuous data, because PCA requires the data to be continuous.

208 7 Experiments

209 7.1 General experiment settings

210 We compare the speed (seconds) and MSE of PCAI with **direct imputation (DI)**, i.e., use an
 211 imputation algorithm directly on the dataset. The imputation approaches used for comparison:
 212 softImpute [20, 40], MissForest [10]¹ and Multiple Imputation by Chained Equation (MICE) [5, 41],
 213 kNN Imputation (KNNI), GAIN [6] are implemented with default configurations. The codes will be
 214 available upon the acceptance of the paper. For PIC, we compare the five fold cross-validation (CV)
 215 score (accuracy, speed) of PIC when dimension reduction is applied on the imputed missing part
 216 (**PIC-reduce**), when dimension reduction is not applied on the imputed missing part (**PIC**), and when
 217 PCA is applied to the imputed version on the full missing data (**DI-reduce**), and when no dimension
 218 reduction is applied to imputed data after direct imputation (**DI**). Here, the default PCA formulation
 219 is PCA-form1, unless specified otherwise. For all PCA computation, the number of eigenvectors is
 220 chosen so that the minimum amount of variance explained is 95%.

221 Details of the datasets used in the experiments are available in Table 1. All experiments are run on
 222 an AMD Ryzen 7 3700X CPU with 8 Cores, 16 processing threads, 3.6GHz, and 16GB RAM. We
 223 terminate an experiment if no result is produced after 6,500 seconds of running or if there arises a
 224 memory allocating issue, and we denote this as **NA** in the result tables.

225 7.2 Performance of PCAI and PIC when the missing values in \mathcal{M} are randomly simulated

Table 2: (MSE, speed) for PCAI and direct imputation (DI) on the Parkinson dataset with $q = 700$.

Imputer	Strategy	missing rate		
		20%	40%	60%
softImpute	PCAI	(0.073, 0.860)	(0.185, 0.774)	(0.305, 0.875)
	DI	(0.072, 4.097)	(0.188, 4.043)	(0.308, 4.467)
MICE	PCAI	(0.091, 139.811)	(0.186, 85.241)	(0.369, 109.815)
	DI	NA	NA	NA
GAIN	PCAI	(0.254, 45.046)	(0.538, 43.938)	(0.779, 43.956)
	DI	(0.608, 69.839)	(1.097, 70.548)	(1.369, 70.293)
missForest	PCAI	(0.064, 188.324)	(0.163, 178.849)	(0.292, 138.085)
	DI	(0.058, 905.002)	(0.160, 692.150)	(0.258, 449.415)
KNNI	PCAI	(0.127, 0.355)	(0.299, 0.398)	(0.466, 0.416)
	DI	(0.113, 0.310)	(0.274, 0.337)	(0.426, 0.372)

226 Note that any datasets can be rearranged so that the first q features are not missing and the remaining
 227 ones are missing. Therefore, without loss of generality, we assume that the first q features of each

¹<https://pypi.org/project/missingpy/>

Table 3: (MSE, speed) for PCAI and DI on the Fashion MNIST dataset with $q = 700$. MissForest results all are NA, and therefore are removed from the tables.

Imputer	Strategy	missing rate		
		20%	40%	60%
softImpute	PCAI	(0.032, 22.408)	(0.066, 22.797)	(0.109, 25.603)
	DI	(0.032, 67.627)	(0.064, 69.349)	(0.107, 77.233)
MICE	PCAI	(0.027, 2218.864)	(0.055, 1374.558)	(0.095, 1641.962)
	DI	NA	NA	NA
GAIN	PCAI	(0.053, 65.730)	(0.091, 68.752)	(0.137, 69.743)
	DI	(0.041, 97.898)	(0.079, 99.049)	(0.125, 96.317)
KNNI	PCAI	(0.055, 1607.850)	(0.115, 2033.153)	(0.180, 2272.370)
	DI	(0.049, 3042.752)	(0.102, 3659.300)	(0.161, 3959.832)

228 dataset are not missing, and the remaining ones contain missing value(s). Then, we simulated missing
 229 data randomly on the missing partition \mathcal{M} with missing rates 20%, 40%, and 60%. Here, a missing
 230 rate of 20% means that 20% of the entries in the missing partition \mathcal{M} are missing. The results for
 231 such experiments are reported in Tables 2, 3, 4. Due to space limit, the results related to PIC on
 Fashion MNIST are reported in the Appendix.

Table 4: Five fold CV results (accuracy, speed) of SVM on Parkinson with $q = 700$.

Imputer	Strategy	missing rate		
		20%	40%	60%
softImpute	PIC-reduce	(0.862, 1.026)	(0.862, 1.137)	(0.862, 1.161)
	PIC	(0.858, 1.008)	(0.858, 1.079)	(0.859, 1.112)
	DI-reduce	(0.861, 4.116)	(0.862, 4.424)	(0.861, 4.718)
	DI	(0.858, 3.775)	(0.858, 3.912)	(0.855, 4.248)
MICE	PIC-reduce	(0.859, 204.605)	(0.861, 256.340)	(0.861, 240.211)
	PIC	(0.858, 524.739)	(0.859, 694.667)	(0.859, 925.426)
	DI-reduce	NA	NA	NA
	DI	NA	NA	NA
GAIN	PIC-reduce	(0.857, 91.086)	(0.852, 102.861)	(0.848, 122.349)
	PIC	(0.851, 89.984)	(0.853, 104.773)	(0.853, 123.233)
	DI-reduce	(0.855, 130.349)	(0.851, 149.864)	(0.851, 181.135)
	DI	(0.846, 129.702)	(0.849, 152.031)	(0.852, 183.67)
missForest	PIC-reduce	(0.859, 204.850)	(0.861, 276.537)	(0.858, 153.783)
	PIC	(0.858, 202.939)	(0.861, 277.067)	(0.858, 153.463)
	DI-reduce	(0.861, 656.948)	(0.862, 729.872)	(0.861, 472.230)
	DI	(0.858, 655.750)	(0.861, 730.013)	(0.858, 472.388)
KNNI	PIC-reduce	(0.858, 0.533)	(0.861, 0.462)	(0.862, 0.625)
	PIC	(0.858, 0.513)	(0.861, 0.462)	(0.862, 0.607)
	DI-reduce	(0.862, 0.696)	(0.862, 0.642)	(0.859, 0.803)
	DI	(0.859, 0.438)	(0.859, 0.45)	(0.858, 0.552)

232

233 From the tables, it is clear that the proposed frameworks reduce the imputation time significantly
 234 while maintaining competitive MSE/classification accuracy compared to DI, in most of the cases.
 235 For example, at the missing rate 20% on the Parkinson dataset (Table 4), when using GAIN for
 236 imputation, the running time of PIC-reduce(91.086s) is much lower compared to DI-reduce (130.349),
 237 the running time of PIC (89.984s) is also much lower compared to DI (129.702). Another example
 238 can be seen from Table 2, for the Parkinson dataset, at 20% missing rate, when PCAI is applied to
 239 missForest, the running time reduces to 188.324s, which is almost 1/5 of the DI (905.002s). Next,

240 on Fashion MNIST (Table 3), it is worth noticing that for MICE, DI cannot give the results due to
 241 memory issue but PCAI can alleviate this issue and deliver the results.

242 For KNNI, the running time for KNNI between the PCAI approach and direct imputation for
 243 Parkinson (Table 2) is not much different. However, for the Fashion MNIST dataset, KNNI using the
 244 PCAI framework obviously deliver a competitive result in a significantly shorter time. Specifically,
 245 KNNI at a missing rate of 20% on Fashion MNIST gives a result after only 1607.850 seconds, while
 246 DI takes up to 3,042.752 seconds. This is because Fashion MNIST (70000 samples) has much more
 247 samples than Parkinson (756 samples), and KNN need to do a lot of pairwise comparison. Therefore,
 248 PCAI and PIC would be extremely helpful for KNNI when the sample size and the number of fully
 249 observed features is large. Note that it does not require the number of features with missing data to
 250 be large or small.

251 From Table 2, we can see that PCAI generates a lot of improvements in MSE for GAIN, in addition
 252 to improvements in speed. This is possibly because PCA reduces the number of features while the
 253 sample size remains the same, making such a deep learning approach more applicable to the newly
 254 reduced data.

255 7.3 Performance on nonrandomly missing data

256 In many fields, the data are missing in a monotone pattern rather than random [44]. Therefore, we
 257 generate one-step monotone missing data on Fashion MNIST by first, randomly choose 20%, 40%,
 258 60% of the samples. Then, we make them become missing by deleting the lower right corner by
 259 deleting the intersection between the last 8 rows and the last 13 columns of each image array. The
 260 results are reported in Table 5. From the table, we can see that PIC-reduce is a great improvement in
 261 speed compared to DI-reduce, and PIC is a significant improvement in speed compared to DI. This
 262 illustrates that PIC can work effectively even for non-randomly missing data.

Table 5: Five fold CV results (accuracy, speed) of SVM on monotone data generated on Fashion MNIST.

Imputer	Strategy	missing rate		
		20%	40%	60%
softImpute	PIC-reduce	(0.889, 409.676)	(0.889, 421.671)	(0.889, 369.49)
	PIC	(0.889, 507.626)	(0.89, 543.309)	(0.889, 480.452)
	DI-reduce	(0.89, 439.268)	(0.89, 528.892)	(0.889, 395.646)
	DI	(0.891, 738.494)	(0.891, 872.616)	(0.89, 646.781)
GAIN	PIC-reduce	(0.886, 462.478)	(0.883, 429.173)	(0.882, 449.786)
	PIC	(0.886, 493.399)	(0.882, 484.232)	(0.881, 496.066)
	DI-reduce	(0.891, 543.803)	(0.89, 431.947)	(0.89, 454.981)
	DI	(0.892, 902.049)	(0.891, 754.794)	(0.891, 780.686)

263 7.4 PIC under different PCA formulations and number of missing features

264 The missing data in these experiments are generated at random as in Section 7.2 and the five fold
 265 cross validation results of SVM on the Gene dataset with $q = 15000, 20000$, are shown in Table 6
 266 and Table 7. From these tables, one can see clearly that for datasets where the number of features
 267 are significantly higher than the number of samples such as Gene, PCA-form2, which is based on
 268 the input data (\mathcal{F} specifically) gives much faster computations compared to PCA-form1, and also is
 269 faster than direct imputation-classification without PCA. In addition, when PCA-form1 is used, even
 270 though PIC and PIC-reduce are faster than PCA on directly imputed data (DI-reduce), they are still
 271 much slower than direct imputation - classification without PCA.

272 Interestingly, the accuracy PIC and PIC-reduce are almost identical to PCA on directly imputed data,
 273 and are higher than direct imputation - classification without PCA. Next, note that the main idea of
 274 the proposed methods is to reduce the dimension of the \mathcal{F} to speed up the imputation. Therefore, we
 275 have made no assumption about the number of features in the missing portion \mathcal{M} . In Table 6 and
 276 Table 7, $q = 15000, 20000$, which means 5,531 and 531 missing features in \mathcal{M} , respectively. This
 277 implies PIC can handle datasets where \mathcal{M} has many features.

Table 6: Five fold CV results (accuracy, speed) of SVM for softImpute based strategies on the Gene dataset when $q = 15000$.

		missing rate		
Strategy		20%	40%	60%
PCA-form1	PIC-reduce	(0.994, 2250.451)	(0.992, 2412.082)	(0.992, 2415.434)
	PIC	(0.992, 2429.114)	(0.992, 2276.354)	(0.992, 2284.414)
	DI-reduce	(0.994, 5018.368)	(0.994, 4529.766)	(0.994, 3785.947)
PCA-form2	PIC-reduce	(0.995, 69.444)	(0.992, 76.393)	(0.992, 85.2)
	PIC	(0.992, 61.451)	(0.992, 68.571)	(0.992, 77.528)
	DI-reduce	(0.995, 80.823)	(0.992, 92.265)	(0.994, 100.751)
No PCA	DI	(0.985, 71.884)	(0.985, 74.812)	(0.985, 92.309)

Table 7: Five fold CV results (accuracy, speed) of SVM for softImpute based strategies on the Gene dataset when $q = 20000$.

		missing rate		
Strategy		20%	40%	60%
PCA-form1	PIC-reduce	(0.994, 2578.910)	(0.994, 4001.717)	(0.994, 3848.950)
	PIC	(0.994, 2583.717)	(0.994, 4144.157)	(0.994, 4057.188)
	DI-reduce	(0.995, 2891.994)	(0.994, 4476.563)	(0.995, 4332.869)
PCA-form2	PIC-reduce	(0.995, 67.753)	(0.992, 73.884)	(0.995, 81.27)
	PIC	(0.995, 59.815)	(0.995, 66.096)	(0.995, 73.079)
	DI-reduce	(0.995, 81.07)	(0.995, 82.407)	(0.995, 91.638)
No PCA	DI	(0.985, 74.06)	(0.985, 71.6)	(0.985, 84.963)

278 8 Conclusion and Remarks

279 We have presented two novel frameworks for datasets where many continuous features are fully
 280 observed, PCAI and PIC, that can speed up imputation algorithms significantly while having com-
 281 petitive accuracy MSE/accuracy compared to direct imputation and alleviate the memory issue for
 282 some imputation approaches such as MICE, kNN. In addition, the frameworks can be used even
 283 when some or all of the missing features are categorical or when the number of missing features
 284 is large. Note that when the sample size is significantly larger than the number of fully observed
 285 features, PCA-form1 should be used since, in such a case, the covariance matrix is much smaller than
 286 \mathcal{F} , making it faster than PCA-form2. On the other hand, when the number of fully observed features
 287 is significantly larger than the sample size, PCA-form2 should be preferred, as the covariance matrix
 288 is bigger than \mathcal{F} itself in such a case. A limitation of the proposed framework is that if there are not
 289 many fully observed continuous features, then due to the computational cost of PCA, the proposed
 290 frameworks may not lead to any improvement in speed.

291 Even though PIC is only introduced for classification, the same strategy can be applied to a regression
 292 problem. We would like to explore that in the future. Moreover, since various dimension reduction
 293 techniques such as sparse PCA [45], incremental PCA [46], truncated SVD [47] have been developed
 294 to suit different scenarios, it is worth investigating different dimension reduction techniques for PCAI
 295 and PIC. In addition, it would be interesting to explore if applying a PCA variant to the missing
 296 partition \mathcal{M} would result in even a more efficient method for datasets with continuous features in the
 297 missing partition.

298 References

299 [1] Zachary C Lipton, David C Kale, Randall Wetzel, et al. Modeling missing data in clinical time
 300 series with rnns. *Machine Learning for Healthcare*, 56, 2016.

- 301 [2] Zhengping Che, Sanjay Purushotham, Kyunghyun Cho, David Sontag, and Yan Liu. Recurrent
302 neural networks for multivariate time series with missing values. *Scientific reports*, 8(1):1–12,
303 2018.
- 304 [3] Mostafa Mehdipour Ghazi, Mads Nielsen, Akshay Pai, M Jorge Cardoso, Marc Modat, Sebastien
305 Ourselin, and Lauge Sørensen. Robust training of recurrent neural networks to handle missing
306 data for disease progression modeling. *arXiv preprint arXiv:1808.05500*, 2018.
- 307 [4] Thu Nguyen, Duy H.M. Nguyen, Huy Nguyen, Binh T. Nguyen, and Bruce A. Wade. Epem:
308 Efficient parameter estimation for multiple class monotone missing data. *Information Sciences*,
309 567:1–22, 2021.
- 310 [5] S van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained
311 equations in r. *Journal of statistical software*, pages 1–68, 2010.
- 312 [6] Jinsung Yoon, James Jordon, and Mihaela Schaar. Gain: Missing data imputation using
313 generative adversarial nets. In *International conference on machine learning*, pages 5689–5698.
314 PMLR, 2018.
- 315 [7] Ayush Garg, Deepika Naryani, Garvit Aggarwal, and Swati Aggarwal. D1-gsa: a deep learning
316 metaheuristic approach to missing data imputation. In *International Conference on Sensing and*
317 *Imaging*, pages 513–521. Springer, (2018).
- 318 [8] Thu Nguyen, Khoi Minh Nguyen-Duy, Duy Ho Minh Nguyen, Binh T. Nguyen, and Bruce Alan
319 Wade. Dper: Direct parameter estimation for randomly missing data. *Knowledge-Based Systems*,
320 240:108082, 2022.
- 321 [9] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for
322 benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- 323 [10] Daniel J Stekhoven and Peter Bühlmann. Missforest-non-parametric missing value imputation
324 for mixed-type data. *Bioinformatics*, 28(1):112–118, 2012.
- 325 [11] Jundong Li, Kewei Cheng, Suhang Wang, Fred Morstatter, Robert P Trevino, Jiliang Tang, and
326 Huan Liu. Feature selection: A data perspective. *ACM computing surveys (CSUR)*, 50(6):1–45,
327 2017.
- 328 [12] Isabelle Guyon, Jiwen Li, Theodor Mader, Patrick A Pletscher, Georg Schneider, and Markus
329 Uhr. Competitive baseline methods set new standards for the nips 2003 feature selection
330 benchmark. *Pattern recognition letters*, 28(12):1438–1444, 2007.
- 331 [13] Bjørn Grung and Rolf Manne. Missing values in principal component analysis. *Chemometrics*
332 *and Intelligent Laboratory Systems*, 42(1-2):125–139, 1998.
- 333 [14] Alexander Ilin and Tapani Raiko. Practical approaches to principal component analysis in the
334 presence of missing values. *The Journal of Machine Learning Research*, 11:1957–2000, 2010.
- 335 [15] Lovedeep Gondara and Ke Wang. Multiple imputation using deep denoising autoencoders.
336 *arXiv preprint arXiv:1705.02737*, 2017.
- 337 [16] Vincent Audigier, François Husson, and Julie Josse. Multiple imputation for continuous
338 variables using a bayesian principal component analysis. *Journal of statistical computation and*
339 *simulation*, 86(11):2140–2156, 2016.
- 340 [17] Dušan Sovilj, Emil Eirola, Yoan Miche, Kaj-Mikael Björk, Rui Nian, Anton Akusok, and
341 Amaury Lendasse. Extreme learning machine for missing data using multiple imputations.
342 *Neurocomputing*, 174:220–231, 2016.
- 343 [18] Karthika Mohan and Judea Pearl. Graphical models for processing missing data. *Journal of the*
344 *American Statistical Association*, pages 1–42, 2021.
- 345 [19] Trevor Hastie, Rahul Mazumder, Jason D Lee, and Reza Zadeh. Matrix completion and
346 low-rank svd via fast alternating least squares. *The Journal of Machine Learning Research*,
347 16(1):3367–3402, 2015.

- 348 [20] Rahul Mazumder, Trevor Hastie, and Robert Tibshirani. Spectral regularization algorithms for
349 learning large incomplete matrices. *Journal of machine learning research*, 11(Aug):2287–2322,
350 2010.
- 351 [21] Sanaz Nikfalazar, Chung-Hsing Yeh, Susan Bedingfield, and Hadi A Khorshidi. Missing data
352 imputation using decision trees and fuzzy clustering with iterative learning. *Knowledge and
353 Information Systems*, 62(6):2419–2437, 2020.
- 354 [22] Md Geaur Rahman and Md Zahidul Islam. Missing value imputation using decision trees and
355 decision forests by splitting and merging records: Two novel techniques. *Knowledge-Based
356 Systems*, 53:51–65, 2013.
- 357 [23] Lane F Burgette and Jerome P Reiter. Multiple imputation for missing data via sequential
358 regression trees. *American journal of epidemiology*, 172(9):1070–1076, 2010.
- 359 [24] Shahidul Islam Khan and Abu Sayed Md Latiful Hoque. Sice: an improved missing data
360 imputation technique. *Journal of big data*, 7(1):1–21, 2020.
- 361 [25] Md Geaur Rahman and Md Zahidul Islam. Missing value imputation using a fuzzy clustering-
362 based em approach. *Knowledge and Information Systems*, 46(2):389–422, 2016.
- 363 [26] Jared S Murray and Jerome P Reiter. Multiple imputation of missing categorical and continuous
364 values via bayesian mixture models with local dependence. *Journal of the American Statistical
365 Association*, 111(516):1466–1479, 2016.
- 366 [27] Collins Leke and Tshilidzi Marwala. Missing data estimation in high-dimensional datasets:
367 A swarm intelligence-deep neural network approach. In *International Conference on Swarm
368 Intelligence*, pages 259–270. Springer, (2016).
- 369 [28] Katarzyna Woźnica and Przemysław Biecek. Does imputation matter? benchmark for predictive
370 models. *arXiv preprint arXiv:2007.02837*, 2020.
- 371 [29] Julie Josse, François Husson, et al. Multiple imputation in principal component analysis.
372 *Advances in data analysis and classification*, 5(3):231–246, 2011.
- 373 [30] Trevor Hastie, Robert Tibshirani, Jerome H Friedman, and Jerome H Friedman. *The elements
374 of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
- 375 [31] Li Qu, Li Li, Yi Zhang, and Jianming Hu. Ppca-based missing data imputation for traffic flow
376 volume: A systematical approach. *IEEE Transactions on intelligent transportation systems*,
377 10(3):512–522, 2009.
- 378 [32] Alfonso Iodice D’Enza, Francesco Palumbo, and Angelos Markos. Single imputation via
379 chunk-wise pca. In *Conference of the International Federation of Classification Societies*, pages
380 75–82. Springer, 2019.
- 381 [33] Abel Folch-Fortuny, Francisco Arteaga, and Alberto Ferrer. Assessment of maximum likelihood
382 pca missing data imputation. *Journal of Chemometrics*, 30(7):386–393, 2016.
- 383 [34] Vincent Audigier, François Husson, and Julie Josse. A principal component method to impute
384 missing values for mixed data. *Advances in Data Analysis and Classification*, 10(1):5–26, 2016.
- 385 [35] Robert Ernest Dear. *A principal-component missing-data method for multiple regression models*.
386 System Development Corporation, 1959.
- 387 [36] Darren T Andrews and Peter D Wentzell. Applications of maximum likelihood principal
388 component analysis: incomplete data sets and calibration transfer. *Analytica Chimica Acta*,
389 350(3):341–352, 1997.
- 390 [37] Sam Roweis. Em algorithms for pca and spca. *Advances in neural information processing
391 systems*, 10, 1997.
- 392 [38] Abel Folch-Fortuny, Francisco Arteaga, and Alberto Ferrer. Pca model building with missing
393 data: New proposals and a comparative study. *Chemometrics and Intelligent Laboratory
394 Systems*, 146:77–88, 2015.

- 395 [39] Aude Sportisse, Claire Boyer, and Julie Josse. Estimation and imputation in probabilistic
 396 principal component analysis with missing not at random data. *Advances in Neural Information*
 397 *Processing Systems*, 33:7067–7077, 2020.
- 398 [40] Alex Rubinsteyn and Sergey Feldman. fancyimpute: An imputation library for python. <https://github.com/iskandr/fancyimpute>, 2016.
- 400 [41] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel,
 401 P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher,
 402 M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine*
 403 *Learning Research*, 12:2825–2830, 2011.
- 404 [42] C Okan Sakar, Gorkem Serbes, Aysegul Gunduz, Hunkar C Tunc, Hatice Nizam, Betul Erdogdu
 405 Sakar, Melih Tutuncu, Tarkan Aydin, M Erdem Isenkul, and Hulya Apaydin. A comparative
 406 analysis of speech signal processing algorithms for parkinson’s disease classification and the
 407 use of the tunable q-factor wavelet transform. *Applied Soft Computing*, 74:255–263, 2019.
- 408 [43] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017.
- 409 [44] Jiang Li, Xiaowei S Yan, Durgesh Chaudhary, Venkatesh Avula, Satish Mudiganti, Hannah
 410 Husby, Shima Shahjouei, Ardavan Afshar, Walter F Stewart, Mohammed Yeasin, et al. Im-
 411 putation of missing values for electronic health record laboratory data. *NPJ digital medicine*,
 412 4(1):1–14, 2021.
- 413 [45] Rodolphe Jenatton, Guillaume Obozinski, and Francis Bach. Structured sparse principal
 414 component analysis. In *Proceedings of the Thirteenth International Conference on Artificial*
 415 *Intelligence and Statistics*, pages 366–373. JMLR Workshop and Conference Proceedings, 2010.
- 416 [46] David A Ross, Jongwoo Lim, Ruei-Sung Lin, and Ming-Hsuan Yang. Incremental learning for
 417 robust visual tracking. *International journal of computer vision*, 77(1):125–141, 2008.
- 418 [47] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness:
 419 Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*,
 420 53(2):217–288, 2011.

421 Checklist

- 422 1. For all authors...
- 423 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 424 contributions and scope? [Yes] . See Table 2, Table 3, Table ??, Table ?? and Section
 425 7.
- 426 (b) Have you read the ethics review guidelines and ensured that your paper conforms to
 427 them? [Yes] .
- 428 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 429 (d) Did you describe the limitations of your work? [Yes] . See the last paragraph of Section
 430 4.
- 431 2. If you are including theoretical results...
- 432 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 433 (b) Did you include complete proofs of all theoretical results? [N/A]
- 434 3. If you ran experiments...
- 435 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 436 mental results (either in the supplemental material or as a URL)? [Yes] . The code will
 437 be available in the form of a Github repository upon acceptance, to preserve anonymity.
- 438 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 439 were chosen)? [Yes] . See Section 7.
- 440 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 441 ments multiple times)? [N/A]

- 442 (d) Did you include the total amount of compute and the type of resources used (e.g., type
443 of GPUs, internal cluster, or cloud provider)? [Yes] . See Section 7.
- 444 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 445 (a) If your work uses existing assets, did you cite the creators? [Yes] . See References list.
- 446 (b) Did you mention the license of the assets? [N/A]
- 447 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
448 . The code can be found on a github repository upon the acceptance of the paper.
- 449 (d) Did you discuss whether and how consent was obtained from people whose data you're
450 using/curating? [N/A]
- 451 (e) Did you discuss whether the data you are using/curating contains personally identifiable
452 information or offensive content? [N/A]
- 453 5. If you used crowdsourcing or conducted research with human subjects...
- 454 (a) Did you include the full text of instructions given to participants and screenshots, if
455 applicable? [N/A]
- 456 (b) Did you describe any potential participant risks, with links to Institutional Review
457 Board (IRB) approvals, if applicable? [N/A]
- 458 (c) Did you include the estimated hourly wage paid to participants and the total amount
459 spent on participant compensation? [N/A]