

ORIGINAL RESEARCH ARTICLE

Benchmarking machine learning missing data imputation methods in large-scale mental health survey databases

Preethi Prakash¹, Kelly Street², Shrikanth Narayanan³, Bridget A. Fernandez^{4,5}, Yufeng Shen⁶, and Chang Shu^{7*} ¹Department of Computer Science, Fu Foundation School of Engineering and Applied Science, Columbia University, New York, NY, United States of America²Department of Population and Public Health Sciences, Division of Biostatistics, Keck School of Medicine, University of Southern California, Los Angeles, CA, United States of America³Viterbi School of Engineering, University of Southern California, Los Angeles, CA, United States of America⁴Department of Pediatrics, Division of Medical Genetics, Children's Hospital Los Angeles and The Saban Research Institute, Los Angeles, CA, United States of America⁵Department of Pediatrics, Keck School of Medicine of USC, University of Southern California, Los Angeles, CA, United States of America⁶Departments of Systems Biology and Biomedical Informatics, and JP Sulzberger Columbia Genome Center, Columbia University Irving Medical Center, New York, NY, United States of America⁷Department of Population and Public Health Sciences, Center for Genetic Epidemiology, Division of Epidemiology and Genetics, Keck School of Medicine, University of Southern California, Los Angeles, CA, United States of America***Corresponding author:**Chang Shu
(april.shu@usc.edu)**Citation:** Prakash P, Street K, Narayanan S, Fernandez BA, Shen Y, Shu C. Benchmarking machine learning missing data imputation methods in large-scale mental health survey databases. *Artif Intell Health*. 2025;2(1):81-92. doi: 10.36922/aih.4406**Received:** August 1, 2024**Revised:** September 17, 2024**Accepted:** October 14, 2024**Published Online:** November 7, 2024**Copyright:** © 2024 Author(s). This is an Open-Access article distributed under the terms of the Creative Commons Attribution License, permitting distribution, and reproduction in any medium, provided the original work is properly cited.**Publisher's Note:** AccScience Publishing remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.**Abstract**

Databases tied to mental and behavioral health surveys suffer from the issue of missing data when participants skip the entire survey, which affects the data quality and sample size. These missing data patterns were investigated and the imputation performance was evaluated in Simons Foundations Powering Autism Research for Knowledge, a large-scale autism cohort consists of over 117,000 participants. Four common methods were assessed – Multiple imputation by chained equations (MICE), K-nearest neighbors (KNN), MissForest, and multiple imputation with denoising autoencoders (MIDAS). In a complete subset of 15,196 autism participants, three types of missingness patterns were simulated. We observed that MIDAS and KNN performed the best as the random missingness rate increased and when blockwise missingness was simulated. The average computational times were each 10 min for MIDAS and KNN, 35 min for MissForest, and 290 min for MICE. MIDAS and KNN both provide promising imputation performance in mental and behavioral health survey data that exhibit blockwise missingness patterns.

Keywords: Missing data; Mental health survey; Imputation methods; Machine learning**1. Introduction**

Large-scale biobank databases in mental and behavioral health such as Simons Foundations Powering Autism Research for Knowledge (SPARK), UK Biobank, and All

of Us have empowered researchers to investigate the genetic and environmental risk factors associated with mental and behavioral disorders among more than 100,000 subjects.¹⁻³ Self-reported surveys and questionnaires such as the social communication questionnaire (SCQ),⁴ repetitive behavior scale-revised (RBS-R),⁵ and developmental coordination disorder questionnaire (DCDQ)⁶ are commonly used to quantify mental and behavioral functions at scale. These questionnaires typically consist of a series of related questions and measure responses using ordinal scales with a natural order or rank to indicate the level of agreement known as Likert scales.⁷

However, missingness commonly occurs in the responses to these surveys and questionnaires. The reasons include non-inapplicable or ambiguous questions, and characteristics of the participants themselves including reluctance to answer sensitive questions, incomplete knowledge, and lack of time.⁸ Missingness can also arise at the source level. Specifically, data may have been curated from varying sources with different administered instrument protocols. Certain questions in the survey also may not be relevant to specific demographic groups, such as those that might not apply to young children.

Common types of missing data include missing completely at random (MCAR) and missing not at random (MNAR), with either specific parts of surveys or entire surveys being incomplete.⁹ In MCAR, the probability of missingness is independent of the observed and unobserved data. MAR is a broader class than MCAR in which the missing data is related to the observed but not the unobserved data. On the other hand, the probability of missingness in MNAR data depends on the unobserved missing values. Typically, participants tend to skip entire questionnaires due to unobserved factors, and a form of MNAR missingness referred to as blockwise missingness arises. Blockwise missingness occurs when all responses belonging to the same survey are missing simultaneously for the same participants, forming clustered missing blocks in the overall phenotypic data.

The simplest solution to address blockwise missingness in mental and behavioral questionnaires is to drop participants with missing surveys.¹⁰ However, this option leads to a significant loss of information, reduced sample size, and loss of statistical power when analyzing mental and behavioral questionnaires in biobank data. Another commonly used approach is to impute missing data using statistical and computational methods. Mean, median, and mode substitutions are basic imputation approaches that maintain the original sample size but can lead to biased inferences.¹¹ Specifically, participants who skip certain questionnaires may exhibit different characteristics than those who complete the questionnaires.¹²

More advanced imputation approaches using statistical and computational methods are needed to accurately impute mental and behavioral surveys with blockwise missingness. Here, four commonly used missing data imputation methods were employed – Multivariate imputation by chained equations (MICE), K-nearest neighbors (KNN), non-parametric missing value imputation using random forest (MissForest), and multiple imputation with denoising autoencoders (MIDAS).¹³⁻¹⁶ MICE is one of the most popular methods of multiple imputation originally developed in the early 2000s.¹³ This approach uses a series of regression models to predict each variable with missingness using the remaining variables in the data.¹⁴ KNN is a supervised machine learning algorithm commonly used when the distribution of the data is unknown or difficult to determine.¹⁵ This method performs predictions on the missing data by averaging the K-nearest data points. MissForest is a missing data imputation method based on a random forest developed in 2012. It predicts missing values based on random forest models trained on the complete dataset and imputes missing values iteratively.¹⁶ MIDAS uses a type of unsupervised neural network to predict missing values in the data by reducing the dimensions in the observed data and reconstructing the missing data. MIDAS was recently developed in 2022 and has proven its high accuracy and computational efficiency through systematic tests on simulated and real social science data.¹⁷

Previous studies have not systematically reviewed machine learning-based imputation methods recently developed for the databases tied to mental and behavioral health surveys.¹⁸⁻²² Most psychiatric studies use multiple imputation for handling missing data and have not taken advantage of the latest machine learning-based imputation techniques.¹⁸⁻²² In addition, they have not focused on assessing imputation accuracy in surveys with blockwise missing structures.¹⁸⁻²² This study systematically examines the imputation performance and computational time of these four commonly used missing data imputation methods (MICE, KNN, MissForest, and MIDAS) in the presence of blockwise missingness in mental and behavioral surveys. It uses data from the SPARK, a large-scale autism research study that collects social functioning and behavioral surveys from over 117,000 participants. This study assesses imputation models on both MCAR and MNAR data, identifying the optimal method for each type of missingness pattern. This study conducts a novel exploration of these methods while also addressing the commonly encountered blockwise missingness pattern.

2. Methods

Figure 1 outlines the sample selection and workflow of the study. The four major steps included: (1) preprocessing the

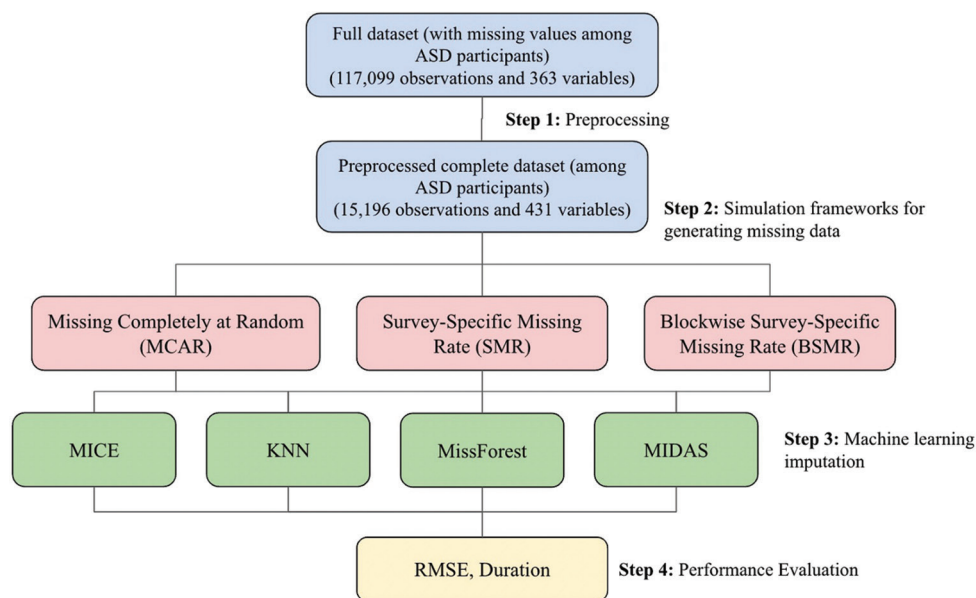


Figure 1. Overview of workflow and study design. (A) The full dataset refers to the original data filtered to only include autism spectrum disorder (ASD) participants. The preprocessed complete dataset refers to the original dataset after filtering to only include ASD participants, dropping incomplete rows, removing variables with extreme rates of missingness, and conducting one-hot-encoding on the categorical variables (which increases the number of variables). (B) Missing completely at random refers to the simulation scenario that randomly converts a specified fraction of the input dataset to missing. Survey-specific missing rate refers to the simulation environment that is tailored to the missingness of the original dataset. Blockwise survey-specific missing rate refers to the simulation environment that is also tailored to the missingness of the original dataset but converts all rows of a given column to missing at once. (C) Multiple imputation by chained equations is an imputation method that employs a series of regression models; MissForest is an imputation method that is based on random forests; Multiple Imputation with Denoising Autoencoders is an imputation method that uses denoising autoencoders; K-nearest neighbors is an imputation method that uses neighboring data points in the feature space. (D) RMSE corresponds to root mean squared error.

data to generate a dataset comprising complete observations, (2) setting up the simulation scenarios for three missing data mechanisms including random missingness, survey-specific missing rates, and blockwise missingness with survey-specific missing rates, (3) conducting the missing data imputation, and (4) evaluating the performance of each model.

2.1. Data source and preprocessing

The dataset used in this study is based on SPARK phenotype V8, consisting of 117,099 participants with autism and 363 variables. It contains information extracted from standardized surveys and parent-reported medical history regarding children with autism. The following eight surveys with <80% missing rates in the full dataset (Table 1) were included in the missing data imputation assessment: individuals registration, basic medical screening, background history, SCQ, RBS-R, DCDQ, Child Behavior Checklist, and area deprivation index.

This dataset was first filtered to remove variables with extreme rates of missingness (~90% or greater), resulting in a drop of 22 variables. The dataset was then modified to remove any rows with missing information. This resulted in 15,196 participants with autism and 347 variables.

Table 1. Percentage of subjects who did not complete each individual survey among all 117,099 participants with autism in SPARK

Survey name	Percentage of subjects who did not complete corresponding survey (%)
Individuals registration	0
Basic medical screening	39.9
Background history	59.3
Area deprivation index	35.1
SCQ	51.3
RBS-R	63.8
DCDQ	72.9
Vineland	82.2
Intelligence quotient	95.3
CBCL	99.6

Note: SCQ: Social communication questionnaire; RBS-R: Repetitive behavior scale-revised; and DCDQ: Developmental coordination disorder questionnaire; are surveys commonly used to quantify the mental and behavioral functions at scale.

Abbreviation: CBCL: Child Behavior Checklist.

One-hot encoding was used to transform the categorical variables in this dataset, resulting in 15,196

participants with autism and 431 variables. The preprocess method from the caret package in R was used to normalize each variable with a mean of 0 and a standard deviation of 1. This was mainly to allow for comparable root mean squared error (RMSE) metrics across all variables that are commonly used in similar studies.^{21,23,24}

This preprocessed complete dataset of participants with autism was used to simulate different missing data mechanisms and assess the accuracy of various imputation methods.

2.2. Three simulation scenarios for missing data mechanisms

Three simulation scenarios were constructed for missing data mechanisms in mental and behavioral surveys as outlined in Figure 2.

2.2.1. MCAR

The first missing data simulation scenario, referred to as MCAR, introduces missingness completely at random by converting a specific percentage of the preprocessed complete dataset to missing. To observe the imputation performance as the missing rate gradually increases, MCAR was implemented with missing rates from 10% to 90% in 10% intervals for all variables in the dataset.

2.2.2. MNAR: SMR

The second missing data simulation scenario is SMR, in which the proportion of missing values in each column is dependent on the survey type that it belongs to. SMR

is tailored to mirror the missing rates in the full SPARK dataset by reusing the same proportions of missing values for each survey (Table 1).

2.2.3. MNAR: BSMR

The last missing data simulation scenario, referred to as BSMR, incorporates blockwise missingness with survey-specific missing rates. Instead of randomly selecting a specific portion of each column to be converted to missing as in SMR, a proportion of participants is randomly selected to have completely missing values for all surveys of a particular survey type. In other words, every column of a specific survey type contains the same missing rows. This resembles real data more closely when subjects skip the entire survey.

2.3. Machine learning imputation

For each missing data simulation scenario described in the previous section, multiple machine learning models were used to impute the missing values. The generated incomplete datasets were passed through the following imputation algorithms to compute the predicted values. A separate set of 10 datasets with 20% randomly selected missing values was used to conduct hyperparameter tuning on each of these models.

2.3.1. MICE

This study used the MICE¹³ (version 3.16.0) package in R which employs a multiple imputation model. It uses a

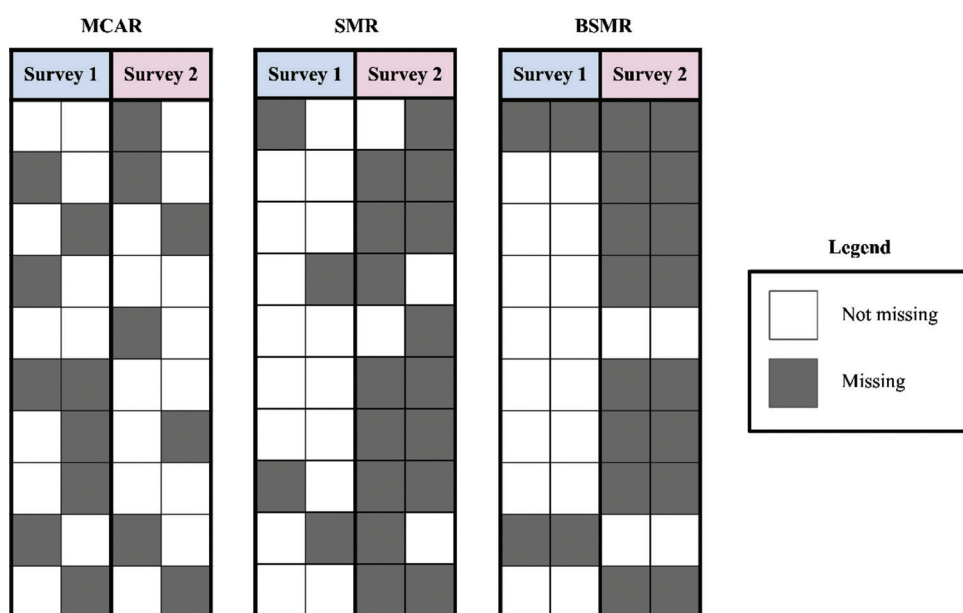


Figure 2. Visualization of the three missing data simulation scenarios explored in this study. On the left is Missing Completely at Random with a 40% missing rate. In the middle is Survey-specific missing rate with a 20% missing rate for Survey 1 and 80% missing rate for Survey 2. On the right is blockwise survey-specific missing rate with a 20% missing rate for Survey 1 and 80% missing rate for Survey 2.

concept called fully conditional specification, in which each incomplete variable is imputed by a different model. It generates multiple imputed datasets that are averaged to retrieve the final imputed data. Since MICEs employ a regression-based approach, hyperparameter tuning was not performed.

2.3.2. KNN

KNNImputer is a method in Python's Scikit-learn package²⁵ (version 0.22) and was used to study the KNN algorithm. KNNImputer predicts each sample's missing values using the average value from the closest data points in the training set. Hyperparameter tuning was used to select the optimal value for the number of nearest neighbors used during imputation.

2.3.3. MissForest

MissForest¹⁶ (version 1.5) is an R package which uses a random forest approach to impute missing values, building multiple decision trees to make predictions using the other remaining features. By averaging several classification or regression trees, MissForest employs out-of-bag error estimates and can capture complex, non-linear relationships. Hyperparameter tuning was used to select the optimal values for the number of trees and the maximum number of iterations.

2.3.4. MIDAS

MIDASpy²⁶ (version 1.3.1) is a Python package that was used to study the MIDAS algorithm. It introduces additional missing values into a given dataset and restores these values using an unsupervised neural network called a denoising autoencoder. Then, the resulting model is used to predict the values of the original missing data. Similar to MICE, MIDASpy generates multiple imputed datasets that are averaged to retrieve the final imputed data. Hyperparameter tuning was used to select the optimal values for the input drop, layer structure, and number of epochs.

2.4. Evaluation of imputation performance

For each missing data simulation scenario, missingness was introduced into the complete dataset 10 different times as 10 separate trials. The values in [Table 1](#) correspond to the percentage of subject IDs in the full dataset (with missing values among participants with autism) who are not present in each specific survey. These missing rates were used when generating the missing datasets for the SMR and BSMR simulation scenarios.

The four models were used to impute the missing data, and these imputed values were compared with the true values in the preprocessed complete dataset. In each imputation trial, the RMSE values were calculated for each

column using the postResample method from the caret package (Version 6.0 – 94) in R. To retrieve the RMSE value for an imputed column, the following formula was used:

$$\text{RMSE} = \sqrt{\frac{\sum (\hat{y}_i - y_i)^2}{n}}$$

Where \hat{y}_i are predicted values and y_i are observed values. As indicated by the equation, the square of the difference between the predicted and observed value was summed across each item in the column that was imputed. This value was then divided by the total number of imputed items and the square root of this value was stored as the column's RMSE.

These column-specific RMSEs were averaged across all columns in the dataset. Then, these RMSEs were again averaged across the 10 trials for each simulation setting. This resulted in a mean overall RMSE for each simulation scenario. These error values were then compared for every simulation scenario between each imputation method.

SCQ summary score, RBS-R summary score, and DCDQ summary score evaluate the social communication function, severity of repetitive behaviors, and motor functions, respectively, in study participants with autism. They were calculated based on corresponding questionnaires. The RMSE values of these specific mental and behavior summary scores were also compared between the four imputation methods across each simulation scenario.

Finally, the total computation time was assessed for the four imputation methods during the BSMR simulation scenario, which was chosen since it is closest in nature to missingness in real survey data.

3. Results

3.1. Overview of full dataset and missingness patterns

The full dataset used in this study consists of 117,099 study participants with autism. Slightly more than half of the participants (51.3%) did not complete SCQ survey, which screens for social functioning; 63.8% did not complete RBS-R survey on repetitive behaviors; and 72.9% did not complete DCDQ survey on motor functions ([Table 1](#)). A total of 34,067 participants have medium missing rates between 20% and 80% among 363 total questions ([Table 2](#)), 37,710 participants exhibit low missing rates (<20%), and 45,322 participants exhibit high missing rates (>80%, [Table 2](#)).

When compared to female participants, there are slightly more male participants with high and low missing

rates. Around 39% of male participants have high missing rates, which is slightly larger than the 37% of female participants, while 33.5% of male participants have low missing rates, and only around 28% of female participants have low missing rates.

For individuals between ages 2 and 18, around 22% of these participants have medium missing rates. The missing rates of these individuals are more concentrated toward extreme values since around 39% have either low or high missing rates or 22% exhibit medium missing rates. For individuals below 2 years of age, around 40% have medium missing rates. Around 62% of individuals above 18 years of age have medium missing rates, whereas nearly 0% exhibit low missing rates.

Close to half of the self-reported white participants, Native Hawaiian participants, and individuals who

Table 2. Demographic characteristics of sample organized by low (<20%), medium (20 – 80%), and high (>80%) missing rate in SPARK

	Missing rate			P-value
	Low missing rate (<20%)	Medium missing rate (20 – 80%)	High missing rate (>80%)	
Number of Subjects	37,710 (32.2)	34,067 (29.1)	45,322 (38.7)	
Sex (%)				<0.001
Male	29,460 (33.5)	24,030 (27.3)	34,412 (39.1)	
Female	8,250 (28.3)	10,037 (34.4)	10,910 (37.4)	
Age (%)				<0.001
<2 years	456 (28.5)	636 (39.7)	509 (31.8)	
2 – 5 years	9,773 (38.0)	6,189 (24.1)	9,726 (37.9)	
6 – 11 years	16,511 (39.1)	9,230 (21.9)	16,463 (39.0)	
12 – 18 years	10,966 (38.4)	6,217 (21.7)	11,401 (39.9)	
>18 years	4 (~0.0)	11,795 (62.0)	7,223 (38.0)	
Race (%)				<0.001
White	28,727 (47.3)	17,968 (30.0)	14,093 (23.2)	
African American	2,063 (37.8)	1,373 (25.2)	2,021 (37.0)	
Asian	876 (35.0)	645 (25.7)	988 (39.4)	
Native American	180 (37.4)	141 (29.3)	160 (33.3)	
Native hawaiian	55 (43.0)	29 (22.7)	44 (34.4)	
Multiple races	4,155 (48.3)	2,203 (25.6)	2,249 (26.1)	
Other	1,654 (4.2)	11,708 (30.0)	25,767 (65.9)	

Note: Proportion of missing variables for each subject was calculated in the full dataset of this study containing 117,099 total participants with autism.

identified as “Multiple Races” have low missing rates. The rates of missingness for self-reported African American, Asian, and Native American individuals are concentrated toward the extreme values, with more than 30% exhibiting high missing rates, while <25% of the participants who were self-identified as White or “Multiple Races” reported high missing rates. Those who self-reported themselves as an “Other” race exhibit large amounts of missingness since around 66% have missing rates larger than 80%.

3.2. Sample characteristics of complete dataset and simulation of three missingness patterns

To assess the imputation performance of the four popular missing data imputation methods (MICE, KNN, MissForest, and MIDAS), a preprocessed complete dataset with 15,196 participants with autism (Table 3, details in

Table 3. Sample characteristics in the preprocessed complete dataset containing 15,196 participants

	Number of observations (percentage) or mean (standard deviation)
Number of subjects	15,196
Sex (%)	
Male	11,901 (78.3)
Female	3,295 (21.7)
Age (%)	
<2 years	61 (0.4)
2 – 5 years	3,029 (19.9)
6 – 11 years	8,442 (55.6)
12 – 18 years	3,664 (24.1)
>18 years	0 (0.0)
Race (%)	
White	11,938 (78.6)
African American	656 (4.3)
Asian	331 (2.2)
Native American	71 (0.5)
Native Hawaiian	22 (0.1)
Multiple races	1,649 (10.9)
Other	529 (3.5)
Summary scores (mean [SD])	
SCQ score	21.72 (7.09)
RBS-R score	35.16 (20.50)
DCDQ score	37.87 (12.73)

Notes: This table includes the number of observations and percentage breakdowns of sex, age, and race as well as means and standard deviations for the summary scores of the; SCQ: Social Communication Questionnaire; RBS-R: Repetitive behavior scale-revised; and DCDQ: Developmental coordination disorder questionnaire.

Methods) was first obtained. Around 78% of participants with complete data are male and 22% are female. The male-to-female ratio is 3.5:1, which aligns with the sex ratio among subjects with autism in the general population. About half of the individuals with complete data are between 6 and 11 years of age. Only 0.4% of subjects are under 2 years of age while none are above 18. About 79% of participants were self-identified as white. The category with the second largest number of participants is “Multiple Races” (10.9%), followed by African American (4.3%), “Other” (3.5%), and Asian (2.2%). The number of participants who are Native American or Native Hawaiian are below 1%. In the preprocessed complete dataset, the SCQ, RBS-R, and DCDQ scores have average values of 21.72, 35.16, and 37.87, respectively.

All variables were standardized with a mean of zero and standard deviation of 1 so that the imputation error, calculated as RMSE, can be interpreted as the average deviation of the predicted scores from the true scores in units of standard deviation. To assess the performance of the missing data imputation methods, missing values were introduced to the preprocessed complete dataset with 15,196 participants with autism. First, to simulate the scenario on MCAR, a random subset of values across the entire dataset was converted to missing values. Ten incomplete datasets were generated for each missingness percentage (10 – 90%). Second, to examine the performance of the imputation methods on MNAR patterns, 10 incomplete datasets were randomly generated for the SMR and BSMR simulation scenarios separately. When doing so, the missing rates in the original SPARK

dataset were used (Table 1) to reflect the missingness distribution present in the real data.

3.3. Performance of imputation on overall dataset

The four imputation methods were applied to the incomplete datasets in each of the three simulation scenarios (Figure 3). The imputed values were compared with the actual values in the complete dataset, and the RMSE values were calculated. RMSE can be interpreted as the average deviation of the predicted scores from the true scores in units of standard deviation since all variables were standardized. Lower RMSE values correspond to higher accuracy in missing value imputation.

In the MCAR scenario, the imputation error for all models generally rose as the missing rate increased. MissForest has the lowest overall RMSE (ranging between 0.73 and 1.0), outperforming the other methods especially when the missing rate was low (Figure 3, left panel). However, as the percentage of missing values increased, the performance of KNN and MIDAS became comparable to that of MissForest. MICE outperformed KNN and MIDAS between 20% and 60% of random missingness but performed considerably worse than all other models for the remaining missing rates.

In the MNAR scenarios, all models exhibited an increase in imputation error in the BSMR scenario when compared to SMR. MissForest produced the lowest error rate in the SMR scenario, with an RMSE of 0.83, but did not perform as well during the BSMR scenario that simulated blockwise missingness. MissForest also exhibited larger variations in

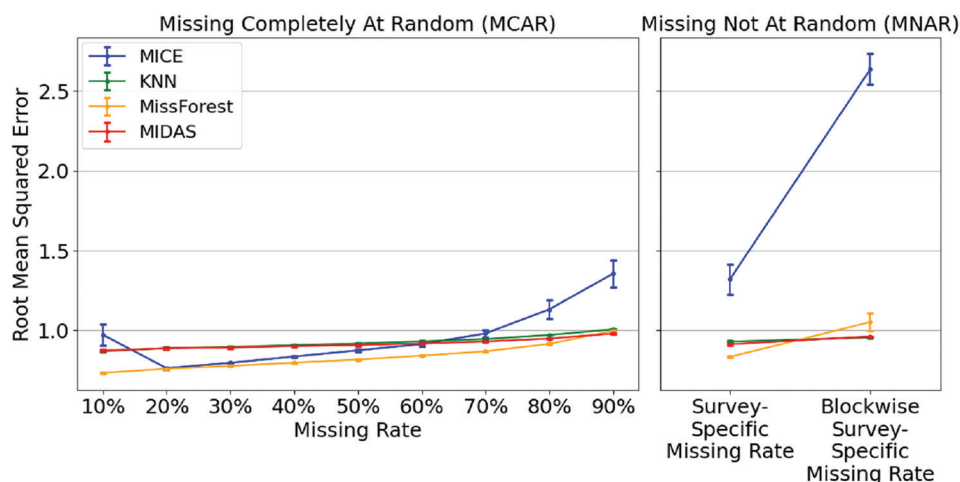


Figure 3. Evaluation of imputation performance based on overall root mean squared error (RMSE). Values across the 10 trials using the missing completely at random simulation scenario (left). Overall RMSE values across the 10 missing not-at-random trials in the survey-specific missing rate and blockwise missingness with Survey-specific missing rate simulation scenarios (right).

Abbreviations: KNN: K-Nearest Neighbors; MICE: Multiple Imputation by Chained Equations; MIDAS: Multiple Imputation with Denoising Autoencoders; MissForest: Non-parametric missing value imputation using Random Forest.

RMSE (standard deviation = 0.056) in the BSMR scenario than in the SMR scenario (standard deviation = 0.0043). For the BSMR scenario, KNN and MIDAS performed the best with an average RMSE of 0.96, outperforming the other methods especially when the missing rate was low (Figure 3, left panel). The variability of the RMSE was also relatively low for both methods, with a standard deviation of 0.0066 for KNN and 1e-6 for MIDAS. MICE performed worse than the other imputation methods in both SMR and BSMR scenarios. Especially in the BSMR scenario, the RMSE value was significantly higher at 2.64 with a relatively large standard deviation of 0.098.

For every simulation scenario, the difference in imputation performance on overall RMSE between KNN and MIDAS was marginal. Both models produced very similar results throughout the experiment and for each simulation scenario besides BSMR, they typically performed slightly worse than MissForest.

3.4. Performance of imputation on mental and behavioral summary scores

For every simulation scenario, the mean and standard deviations of RMSE values for the SCQ, RBS-R, and DCDQ scores were computed across the 10 trials as displayed in Figure 4. The relative performance of the four models was generally consistent across the three summary scores.

In the MCAR scenario, MissForest consistently outperformed KNN and MIDAS when imputing all three summary scores. The MICE model exhibited a steep incline in error as the missing rate was incremented. It performed the best until the missing rate was increased to 50%, after which it was surpassed by the remaining models. MICE is ideal for lower rates of random missingness but begins to perform exponentially worse as the rate gets larger. In fact, the MICE model produced the largest RMSE among the four methods at a 90% missing rate. For missing rates that are 50% and above, MissForest is the ideal model since it had the lowest errors among the four methods.

The MissForest model performed the best in the SMR scenario. However, each method, especially MICE and MissForest, exhibited error rates that rose sharply when the missing values became blocked by survey type in the BSMR scenario. In the BSMR scenario, KNN and MIDAS exhibited the lowest error rates with MissForest performing slightly worse. MICE performed considerably worse than the remaining models in the BSMR scenario.

3.5. Computational time

When comparing the computational times of the four models, the BSMR simulation scenario was used since this environment most closely resembles the missingness

patterns in the real data when participants skip an entire survey in SPARK.

As shown in Figure 5, MIDAS and KNN not only had similar overall error rates but also exhibited comparable imputation times of around 10 – 13 min. MissForest had a median imputation time of slightly <30 min. On the other hand, MICE had a median imputation time of around 285 min, which was significantly larger than those of the remaining models. The difference in computational time between implementations in R and Python is negligible.²⁶

4. Discussion

The establishment of biobank databases has enabled the collection of self-reported mental and behavioral surveys at scale.¹⁻³ SPARK has gathered social and behavioral survey data from about 100,000 individuals¹ and there is ongoing collection of more survey data on existing participants. UK Biobank has measurements on lifetime depressive disorder, cognitive function, attention, and impulsivity from about 150,000 participants.^{2,27,28} All of Us also has strategic plans to collect mental and behavioral surveys at scale.³ However, the data quality and statistical power are compromised by missing data. Recent advances in machine learning methods have inspired novel missing data imputation approaches with increased accuracy and computational efficiency.¹³⁻¹⁶ Previous studies either have not reviewed these newly developed imputation methods or have not focused on assessing imputation accuracy in mental and behavioral surveys that exhibit blockwise missing structures.¹⁸⁻²²

Our study provided insights on the missingness pattern in SPARK, a large-scale cohort with autism, and assessed the imputation accuracy and computational time of four popular missing data imputation methods—MICE, KNN, MissForest, and MIDAS. This was done by simulating three missingness scenarios in mental and behavioral surveys, including SCQ, RBS-R, and DCDQ. We observed that 50 – 70% of participants with autism did not complete SCQ, RBS-R, and DCDQ surveys and the dataset exhibited blockwise missing structures. The missing rates also varied by sex, age, and race. Overall, KNN and MIDAS showed relatively stable performance with increasing missing rate in the MCAR scenario and slightly higher imputation error when blockwise missingness is introduced in the MNAR scenarios. The error rate increased more significantly in MICE and MissForest in both MCAR and MNAR scenarios, with a particularly notable surge in error rate for MICE when blockwise missing structures were introduced. When imputing SCQ, RBS-R, and DCDQ summary scores in the MCAR scenario, MICE had the lowest error rate when the missing rate was low, while MissForest had the lowest

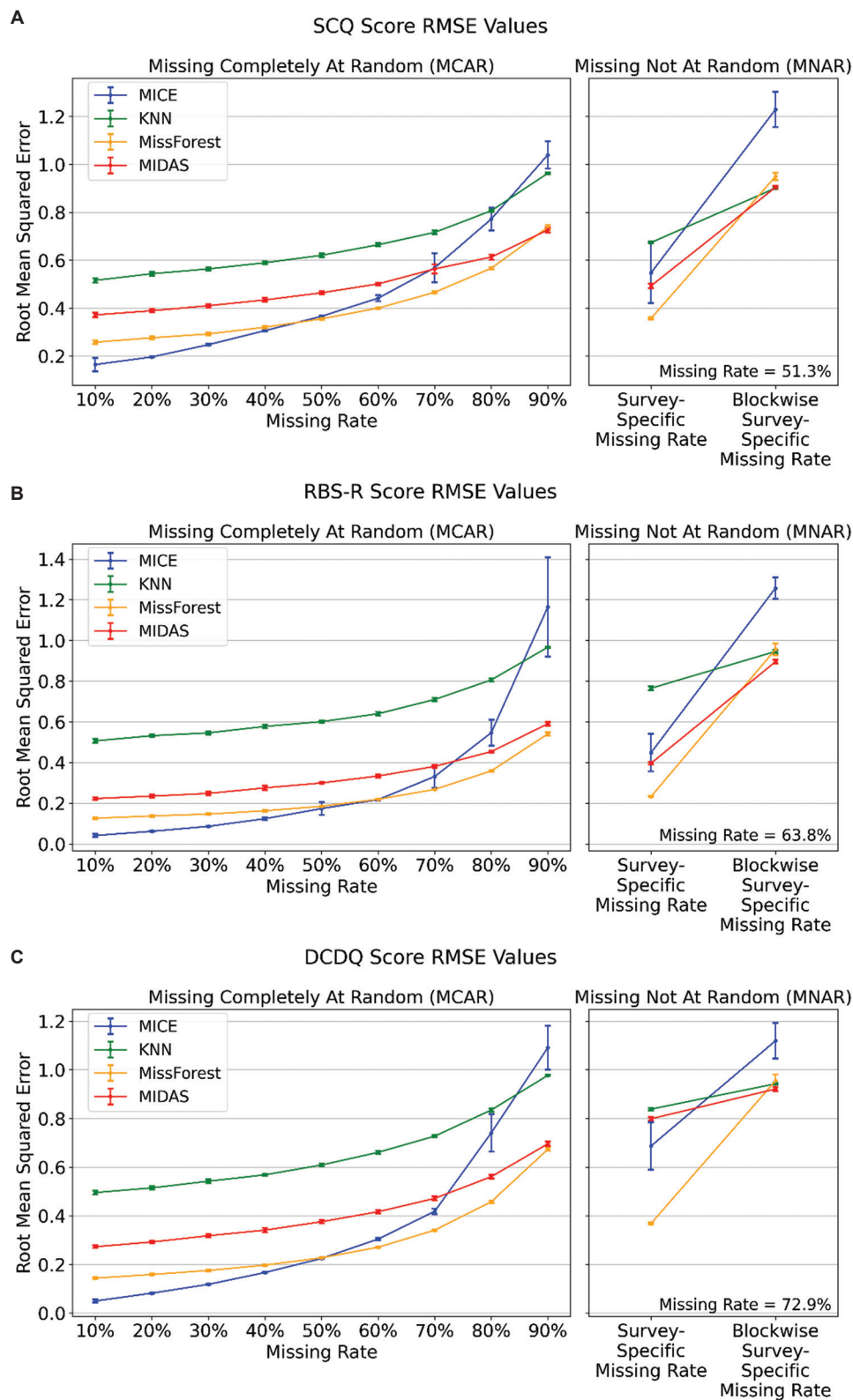


Figure 4. Imputation performance on summary scores from mental health surveys. Root mean squared error (RMSE) values for imputing social communication questionnaire scores (A), Repetitive Behavior Scale-Revised scores (B), Developmental coordination disorder questionnaire scores (C) across the Missing Completely at Random (MCAR) and missing not at random (MNAR) trials. RMSE values for the across the MCAR and MNAR trials.

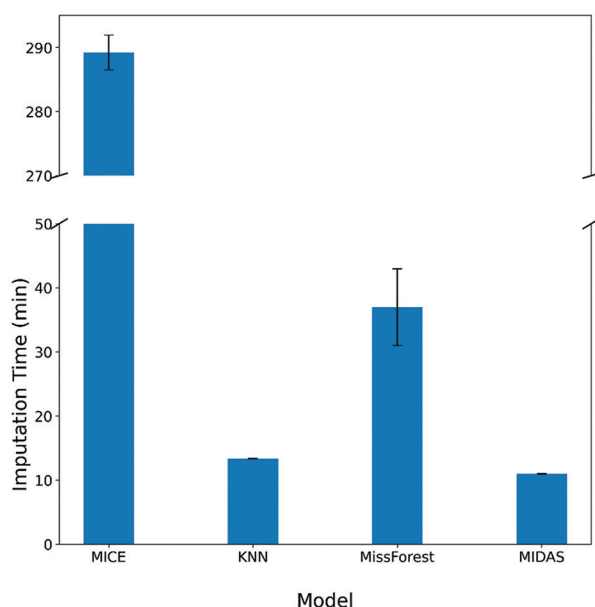


Figure 5. Total imputation times (in minutes) and standard deviations of each model for the 10 trials in the Blockwise Missingness with BSMR scenario. The total sample size is 15,196.

Abbreviations: KNN: K-Nearest Neighbors; MICE: Multiple Imputation by Chained Equations; MIDAS: Multiple Imputation with Denoising Autoencoders; MissForest: Non-parametric missing value imputation using Random Forest; BSMR: Survey-specific missing rate.

error rate when the missing rate was high. However, in the presence of blockwise missingness in the MNAR scenario, MIDAS was consistently the best-performing model across all three summary scores, with KNN and MissForest having similar or slightly higher error rates. The results of this study suggested that some models like MICE are sensitive to high missing rates and blockwise missing structures, while MIDAS and KNN may perform better in the overall dataset and specific summary scores in the presence of blockwise missingness. The average computational times were each 10 min for MIDAS and KNN to impute 15,196 subjects with blockwise missingness, about 35 min for MissForest, and about 290 min for MICE. These results highlight the computational efficiency in machine learning imputation algorithms even in highly complex neural network models in MIDAS. Newly developed imputation models have better optimization in their algorithms and take advantage of parallel computing to reduce the computational time.

Our results show the potential to impute missing data in large-scale databases with mental and behavioral surveys, especially imputing summary scores based on medical history and neurodevelopmental measures. When the data exhibits blockwise missingness, the imputation error increases, but models such as MIDAS and KNN can still provide imputed results that are relatively stable

and accurate. This shows that when a block of correlated variables in one survey is completely missing, other related surveys or medical history can also provide relevant information for imputation. The choice of imputation methods may depend on the overall missing rate and missingness patterns in a dataset.

The strength of our study is that a large-scale collection of mental and behavioral surveys in SPARK was utilized to simulate the missingness patterns, particularly with blockwise missing structures that are commonly observed in mental health databases. This study also systematically assessed the latest missing data imputation approaches like MIDAS. The limitation is that the complete data with missing data simulation primarily comes from adolescents. Despite the inclusion of various racial groups in the simulation, most participants are white. Assessment in other types of large-scale mental and behavioral surveys with adults and minority groups is warranted for future studies.

Missing data imputation is widely used in national surveys with mental and behavioral surveys. For example, the National Survey on Drug Use and Health (NSDUH) has been providing imputation-revised variables by the predictive mean neighborhood methods since 1999.²⁹ There is also the recent phenotype imputation model developed in the UK Biobank, which has shown increased power for genetic studies.³⁰ As biobanks and national surveys collect more large-scale data on mental and behavioral surveys, missing data imputation will produce more accurate imputed values and become an integral part of analysis to maximize the use of the data.

5. Conclusion

Our study underscores the efficacy of advanced imputation techniques, such as MIDAS and KNN, in addressing missing data within large-scale mental and behavioral surveys. Our findings showcase that for similar databases with mental and behavioral surveys on autism, dementia, and other disorders, machine learning-based imputation methods can be leveraged to effectively recover missing information. This study demonstrates that machine learning methods offer increased performance and faster computation times over traditional algorithms. The performance of these advanced imputation techniques demonstrates their potential to optimize analyses and advance research in mental and behavioral disorders.

Acknowledgments

The authors are extremely grateful to the thousands of individuals and families who are participating in the SPARK. The authors also thank the sites, staff, and

volunteers of the SPARK Clinical Site Network and SFARI for their invaluable contributions.

Funding

This work is supported by Southern California Environmental Health Sciences Center pilot grant from NIH/NIEHS, grant number P30ES007048 (Rob McConnell), and The Tobacco-Related Disease Research Program, grant number T32IR5216 (Xuejuan Jiang) and NIH/NIA, grant number 1RF1AG076124-01A1 (Hussein Yassine).

Conflict of interest

The authors declare that they have no conflicts of interest.

Author contributions

Conceptualization: Chang Shu

Formal analysis: Preethi Prakash

Investigation: Preethi Prakash, Chang Shu

Methodology: Preethi Prakash, Kelly Street, Yufeng Shen, Chang Shu

Writing—original draft: Preethi Prakash, Chang Shu

Writing—review & editing: Kelly Street, Shrikanth Narayanan, Bridget A. Fernandez, Yufeng Shen, Chang Shu

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Availability of data

SPARK Phenotype Dataset is accessible through an application at SFARI Base (<https://base.sfari.org>). All software used in this study is publicly available. The code for simulations and analysis can be found at <https://github.com/AprilShuLab/MissingDataImputation>.

Further disclosure

The paper has been uploaded to medRxiv (doi: 10.1101/2024.05.13.24307231).

References

1. Feliciano P, Daniels AM, Snyder LG, *et al.* SPARK: A US cohort of 50,000 families to accelerate autism research. *Neuron*. 2018;97:488-493.
doi: 10.1016/j.neuron.2018.01.015
2. Davis KAS, Coleman JRI, Adams M, *et al.* Mental health in UK Biobank - development, implementation and results from an online questionnaire completed by 157 366 participants: A reanalysis. *BJPsych Open*. 2020;6:e18.
doi: 10.1192/bjo.2019.100
3. Ramirez AH, Sulieman L, Schlueter DJ, *et al.* The all of Us research program: Data quality, utility, and diversity. *Patterns (N Y)*. 2022;3:100570.
doi: 10.1016/j.patter.2022.100570
4. Chesnut SR, Wei T, Barnard-Brak L, Richman DM. A meta-analysis of the social communication questionnaire: Screening for autism spectrum disorder. *Autism*. 2017;21:920-928.
doi: 10.1177/1362361316660065
5. Hooker JL, Dow D, Morgan L, Schatschneider C, Wetherby AM. Psychometric analysis of the repetitive behavior scale-revised using confirmatory factor analysis in children with autism. *Autism Res*. 2019;12:1399-1410.
doi: 10.1002/aur.2159
6. Van Damme T, Vancampfort D, Thoen A, Sanchez CPR, van Biesen D. Evaluation of the Developmental Coordination Questionnaire (DCDQ) as a screening instrument for co-occurring motor problems in children with autism spectrum disorder. *J Autism Dev Disord*. 2022;52:4079-4088.
doi: 10.1007/s10803-021-05285-1
7. Jebb AT, Ng V, Tay L. A review of key likert scale development advances: 1995-2019. *Front Psychol*. 2021;12:637547.
doi: 10.3389/fpsyg.2021.637547
8. Mirzaei A, Carter SR, Patanwala AE, Schneider CR. Missing data in surveys: Key concepts, approaches, and applications. *Res Soc Adm Pharm*. 2022;18:2308-2316.
doi: 10.1016/j.sapharm.2021.03.009
9. Mack C, Su Z, Westreich D. *Managing Missing Data in Patient Registries: Addendum to Registries for Evaluating Patient Outcomes: A User's Guide, Third Edition*. Rockville, MD: Agency for Healthcare Research and Quality (US); 2018.
10. Khan SI, Hoque ASM. SICE: An improved missing data imputation technique. *J Big Data*. 2020;7:37.
doi: 10.1186/s40537-020-00313-w
11. Phiwhorm K, Saikaew C, Leung CK, Polpinit P, Saikaew KR. Adaptive multiple imputations of missing values using the class center. *J Big Data*. 2022;9:52.
doi: 10.1186/s40537-022-00608-0
12. De Goeij MCM, van Diepen M, Jager KJ, Tripepi G, Zoccali C, Dekker FW. Multiple imputation: Dealing with missing data. *Nephrol Dial Transplant*. 2013;28:2415-2420.
doi: 10.1093/ndt/gft221
13. Van Buuren S, Groothuis-Oudshoorn K. Mice: Multivariate imputation by chained equations in R. *J Stat Softw*. 2011;45:1-67.

- doi: 10.18637/jss.v045.i03
14. Azur MJ, Stuart EA, Frangakis C, Leaf PJ. Multiple imputation by chained equations: What is it and how does it work? *Int J Methods Psychiatr Res.* 2011;20:40-49.
doi: 10.1002/mpr.329
15. Taunk K, De S, Verma S, Swetapadma A. A Brief Review of Nearest Neighbor Algorithm for Learning and Classification. In: *2019 International Conference on Intelligent Computing and Control Systems (ICCS)*. IEEE; 2019:1255-1260.
doi: 10.1109/iccs45141.2019.9065747
16. Stekhoven DJ, Bühlmann P. MissForest--non-parametric missing value imputation for mixed-type data. *Bioinformatics.* 2011;28:112-118.
doi: 10.1093/bioinformatics/btr597
17. Lall R, Robinson T. The MIDAS touch: Accurate and scalable missing-data imputation with deep learning. *Polit Anal.* 2022;30:179-196.
doi: 10.1017/pan.2020.49
18. Shrive FM, Stuart H, Quan H, Ghali WA. Dealing with missing data in a multi-question depression scale: A comparison of imputation methods. *BMC Med Res Methodol.* 2006;6:57.
doi: 10.1186/1471-2288-6-57
19. Peyre H, Leplège A, Coste J. Missing data methods for dealing with missing items in quality of life questionnaires. A comparison by simulation of personal mean score, full information maximum likelihood, multiple imputation, and hot deck techniques applied to the SF-36 in the French 2003 decennial health survey. *Qual Life Res.* 2011;20:287-300.
doi: 10.1007/s11136-010-9740-3
20. Emmanuel T, Maupong T, Mpoeleng D, Semong T, Mphago B, Tabona O. A survey on missing data in machine learning. *J Big Data.* 2021;8:140.
doi: 10.1186/s40537-021-00516-9
21. Xu X, Xia L, Zhang Q, Wu S, Wu M, Liu H. The ability of different imputation methods for missing values in mental measurement questionnaires. *BMC Med Res Methodol.* 2020;20:42.
doi: 10.1186/s12874-020-00932-0
22. Croy CD, Novins DK. Methods for addressing missing data in psychiatric and developmental research. *J Am Acad Child Adolesc Psychiatry.* 2005;44:1230-1240.
doi: 10.1097/01.chi.0000181044.06337.6f
23. Lee JH, Huber JC Jr. Evaluation of multiple imputation with large proportions of missing data: How much is too much? *Iran J Public Health.* 2021;50:1372-1380.
doi: 10.18502/ijph.v50i7.6626
24. Petrazzini BO, Naya H, Lopez-Bello F, Vazquez G, Spangenberg L. Evaluation of different approaches for missing data imputation on features associated to genomic data. *BioData Mining.* 2021;14:44.
doi: 10.1186/s13040-021-00274-7
25. Pedregosa F, Varoquaux G, Gramfort A, et al. Scikit-learn: Machine learning in Python. *J Mach Learn Res.* 2011;12:2825-2830.
26. Lall R, Robinson T. Efficient multiple imputation for diverse data in python and R: MIDASpy and rMIDAS. *J Stat Softw.* 2023;107:1-38.
doi: 10.18637/jss.v107.i09
27. Fawns-Ritchie C, Deary IJ. Reliability and validity of the UK Biobank cognitive tests. *PLoS One.* 2020;15:e0231627.
doi: 10.1371/journal.pone.0231627
28. Schwenen LJS, van Rooij D, Shi H, et al. Diet, physical activity, and disinhibition in middle-aged and older adults: A UK biobank study. *Nutrients.* 2021;13:1607.
doi: 10.3390/nu13051607
29. Grau E, Frechtel P, Odom D, Painter D. A Simple Evaluation of the Imputation Procedures Used in NSDUH. In: *Proceedings of the 2004 Joint Statistical Meetings, American Statistical Association, Section on Survey Research Methods, Toronto, Ontario, Canada [CD-ROM]*. Alexandria, VA: American Statistical; 2004.
30. An U, Pazokitoroudi A, Alvarez M, et al. Deep learning-based phenotype imputation on population-scale biobank data increases genetic discoveries. *Nat Genet.* 2023;55:2269-2276.
doi: 10.1038/s41588-023-01558-w