

# Benchmarking Machine Learning Missing Data Imputation Methods in Large-Scale Mental Health Survey Databases

Preethi Prakash<sup>1</sup>, Kelly Street<sup>2</sup>, Shrikanth Narayanan<sup>3</sup>, Bridget A. Fernandez<sup>4,5</sup>, Yufeng Shen<sup>6</sup>, Chang Shu<sup>7</sup>

1. Department of Computer Science, Columbia University, New York, NY, USA
2. Division of Biostatistics, Department of Population and Public Health Sciences, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA
3. Viterbi School of Engineering, University of Southern California, Los Angeles, CA, USA
4. Division of Medical Genetics, Department of Pediatrics, Children's Hospital Los Angeles and The Saban Research Institute, Los Angeles, CA, USA
5. Department of Pediatrics, Keck School of Medicine of USC, University of Southern California, Los Angeles, CA, USA
6. Department of Systems Biology, Department of Biomedical Informatics, and JP Sulzberger Columbia Genome Center, Columbia University Irving Medical Center, New York, NY, USA.
7. Center for Genetic Epidemiology, Division of Epidemiology and Genetics, Department of Population and Public Health Sciences, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA

## Abstract

Databases with mental and behavioral health surveys suffer from missingness when participants skip the entire survey, affecting the data quality and sample size. We investigated the missing data patterns and evaluate the imputation performance in Simons Powering Autism Research (SPARK), 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, we simulated three types of missingness patterns. We observed that MIDAS and KNN performed the best as the rate of random missingness increased and when blockwise missingness was simulated. The average computational times for MIDAS and KNN were 10 minutes, 35 minutes for MissForest, and 290 minutes 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, machine learning

**Author for Correspondence:** Chang Shu ([april.shu@usc.edu](mailto:april.shu@usc.edu))

## 41 **Introduction**

42 Large-scale biobank databases in mental and behavioral health such as Simons Powering Autism  
43 Research for Knowledge (SPARK), UK Biobank and All of Us have empowered researchers to  
44 investigate the genetic and environmental risk factors associated with mental and behavioral  
45 disorders among more than 100,000 subjects<sup>1-3</sup>. Self-reported surveys and questionnaires such as  
46 the Social Communication Questionnaire (SCQ)<sup>4</sup>, Repetitive Behavior Scale-Revised (RBS-R)<sup>5</sup>  
47 and Developmental Coordination Disorder Questionnaire (DCDQ)<sup>6</sup> are commonly used to  
48 quantify mental and behavioral functions at scale. These questionnaires typically consist of a series  
49 of related questions and measure responses using ordinal scales with a natural order or rank to  
50 indicate level of agreement known as Likert scales<sup>7</sup>.

51  
52 However, missingness commonly occurs in the responses to these surveys and questionnaires. The  
53 reasons include non-inapplicable or ambiguous questions, and characteristics of the participants  
54 themselves including reluctance to answer sensitive questions, incomplete knowledge, and lack of  
55 time. Missingness can also arise at the source level. Specifically, data may have been curated from  
56 varying sources with different administered instrument protocols. Certain questions in the survey  
57 also may not be relevant to specific demographic groups, such as those that might not apply to  
58 young children.

59  
60 Common types of missing data include Missing Completely at Random (MCAR) and Missing Not  
61 at Random (MNAR) with either specific parts of surveys or entire surveys being incomplete<sup>8</sup>. In  
62 MCAR, the probability of missingness is independent of the observed and unobserved data. MAR  
63 is a broader class than MCAR in which the missing data is related to the observed but not the  
64 unobserved data. On the other hand, the probability of missingness in MNAR data depends on the  
65 unobserved missing values. Typically, participants tend to skip entire questionnaires due to  
66 unobserved factors, and a form of MNAR missingness referred to as blockwise missingness arises.  
67 Blockwise missingness occurs when all responses belonging to the same survey are missing  
68 simultaneously for the same participants, forming clustered missing blocks in the overall  
69 phenotypic data.

70  
71 The simplest solution to address blockwise missingness in mental and behavioral questionnaires  
72 is to drop participants with missing surveys<sup>9</sup>. However, this option leads to a significant loss of  
73 information, reduced sample size and loss of statistical power when analyzing mental and  
74 behavioral questionnaires in biobank data. Another commonly used approach is to impute missing  
75 data using statistical and computational methods. Mean, median, and mode substitutions are basic  
76 imputation approaches that maintain the original sample size but can lead to biased inferences<sup>10</sup>.  
77 Specifically, participants who skip certain questionnaires may exhibit different characteristics than  
78 those who complete the questionnaires<sup>11</sup>.

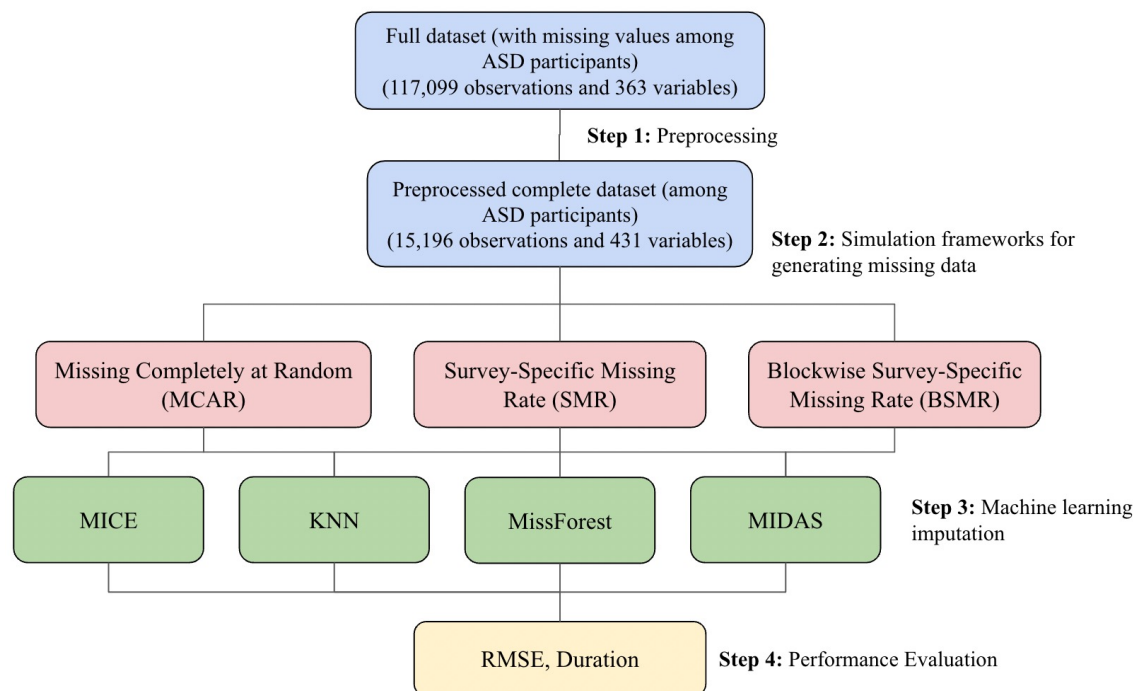
79  
80 More advanced imputation approaches using statistical and computational methods are needed to  
81 accurately impute mental and behavioral surveys with blockwise missingness. Here, we employed  
82 four commonly used missing data imputation methods - Multivariate Imputation by Chained  
83 Equations (MICE), K-Nearest Neighbors (KNN), nonparametric missing value imputation using  
84 Random Forest (MissForest), and Multiple Imputation with Denoising Autoencoders (MIDAS)<sup>12-</sup>  
85 <sup>15</sup>. MICE is one of the most popular methods of multiple imputation originally developed in the  
86 early 2000s<sup>12</sup>. This approach uses a series of regression models to predict each variable with

87 missingness using the remaining variables in the data <sup>13</sup>. KNN is a supervised machine learning  
88 algorithm commonly used when the distribution of the data is unknown or difficult to determine  
89 <sup>14</sup>. This method performs predictions on the missing data by averaging the  $k$  nearest data points.  
90 Nonparametric Missing Value Imputation using Random Forest (MissForest) is a missing data  
91 imputation method based on random forest developed in 2012. It predicts missing values based on  
92 random forest models trained on the complete dataset and imputes missing values iteratively <sup>15</sup>.  
93 Multiple Imputation with Denoising Autoencoders (MIDAS) uses a type of unsupervised neural  
94 network to predict missing values in the data by reducing the dimensions in the observed data and  
95 reconstructing the missing data. MIDAS was recently developed in 2022 and has proven its high  
96 accuracy and computational efficiency through systematic tests on simulated and real social  
97 science data <sup>16</sup>.

98  
99 Previous studies have not systematically reviewed new imputation methods in the databases with  
100 mental and behavioral health surveys<sup>17-21</sup>. Additionally, they have not focused on assessing  
101 imputation accuracy in surveys with blockwise missing structures<sup>17-21</sup>. This study systematically  
102 examines the imputation performance and computational time of these four commonly used  
103 missing data imputation methods (MICE, KNN, MissForest and MIDAS) in the presence of  
104 blockwise missingness in mental and behavioral surveys. It uses data from the Simons Powering  
105 Autism Research for Knowledge (SPARK), a large-scale autism research study that collects social  
106 functioning and behavioral surveys from over 117,000 participants. This study assesses imputation  
107 models on both MCAR and MNAR data, identifying the optimal method for each type of  
108 missingness pattern. This study conducts a novel exploration of these methods while also  
109 addressing the commonly encountered blockwise missingness pattern.

## 110 111 **Methods**

112 **Figure 1** outlines the sample selection and workflow of the study. The four major steps included  
113 (1) preprocessing the data to generate a dataset comprised of complete observations, (2) setting up  
114 the simulation scenarios for three missing data mechanisms including random missingness,  
115 survey-specific missing rates, and blockwise missingness with survey-specific missing rates, (3)  
116 conducting the missing data imputation, and (4) evaluating the performance of each model.



117  
118 **Figure 1. Overview of workflow and study design.** a) The full dataset refers to the original data filtered  
119 to only include ASD participants. The preprocessed complete dataset refers to the original dataset after  
120 filtering to only include ASD participants, dropping incomplete rows, removing variables with extreme  
121 rates of missingness, and conducting one-hot-encoding on the categorical variables (which increases the  
122 number of variables). b) MCAR refers to the simulation scenario which randomly converts a specified  
123 fraction of the input dataset to missing. SMR refers to the simulation environment that is tailored to the  
124 missingness of the original dataset. BSMR refers to the simulation environment that is also tailored to the  
125 missingness of the original dataset, but converts all rows of a given column to missing at once. c) MICE  
126 is an imputation method that employs a series of regression models; MissForest is an imputation method  
127 that is based on random forests; MIDAS is an imputation method that uses denoising autoencoders; KNN  
128 is an imputation method that uses neighboring data points in the feature space. d) RMSE corresponds to  
129 Root Mean Squared Error.

130  
131 **1. Data Source and Preprocessing**  
132 The dataset used in this study is based on SPARK phenotype V8, with 117,099 participants with  
133 autism and 363 variables. It contains information extracted from standardized surveys and parent-  
134 reported medical history regarding children with autism. The following 8 surveys with <80%  
135 missing rates in the full dataset (**Table 1**) were included in missing data imputation assessment:  
136 Individuals Registration, Basic Medical Screening, Background History, Social Communication  
137 Questionnaire (SCQ), Repetitive Behavior Scale-Revised (RBS-R), Developmental Coordination  
138 Disorder Questionnaire (DCDQ), Child Behavior Checklist (CBCL), and Area Deprivation Index

139 (ADI).

Survey Name	Percentage of Subjects Who Did Not Complete Corresponding Survey (%)
Individuals Registration	00.0
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

**Table 1. Percentage of subjects who did not complete each individual survey among all 117,099 ASD participants in SPARK.** Social Communication Questionnaire (SCQ), Repetitive Behavior Scale-Revised (RBS-R), Developmental Coordination Disorder Questionnaire (DCDQ), and Child Behavior Checklist (CBCL) are surveys commonly used to quantify the mental and behavioral functions at scale.

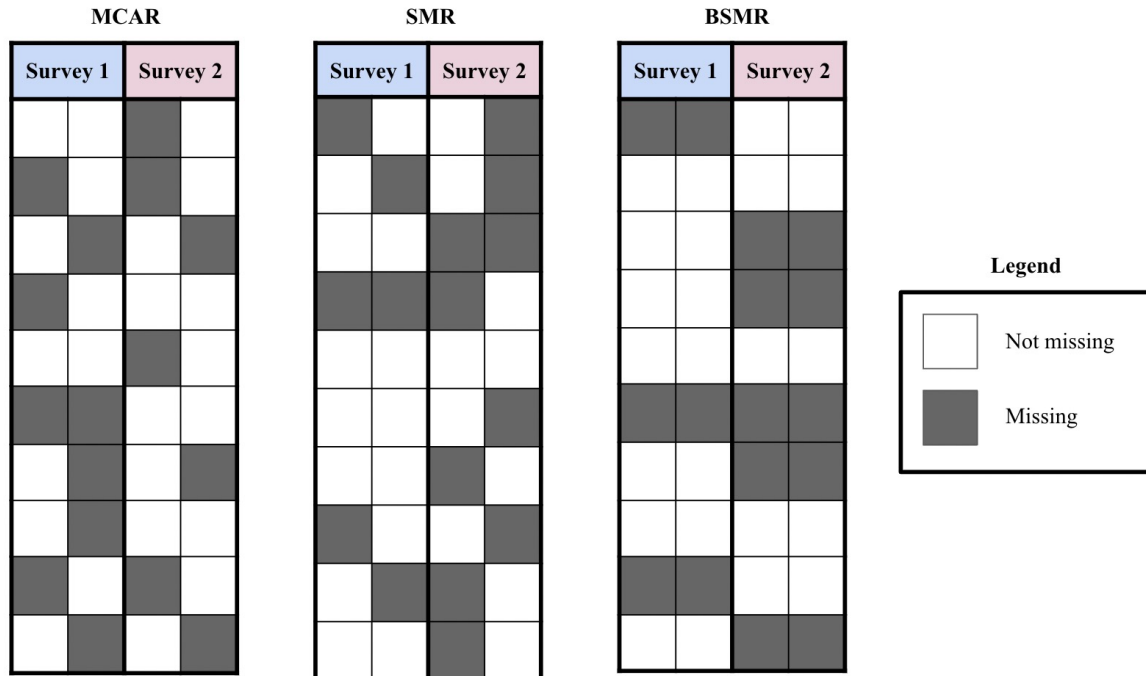
140  
141  
142 This dataset was first filtered to remove variables with extreme rates of missingness (~90% or  
143 greater), resulting in a drop of 22 variables. The dataset was then modified to remove any rows  
144 with missing information. This resulted in 15,196 participants with autism and 347 variables.

145  
146 One-hot encoding was used to transform the categorical variables in this dataset, resulting in  
147 15,196 participants with autism and 431 variables. The `preProcess` method from the `caret`  
148 package in R was used to center and scale each column to have a mean of 0 and standard deviation  
149 of 1. This was mainly to allow for comparable Root Mean Squared Error (RMSE) metrics across  
150 all variables.

151  
152 This preprocessed complete dataset of participants with autism was used to simulate different  
153 missing data mechanisms and assess the accuracy or various imputation methods.

## 154 155 **2. Three Simulation Scenarios for Missing Data Mechanisms**

156 We simulated three simulation scenarios for missing data mechanisms in mental and behavioral  
157 surveys as outlined below and in **Figure 2**.



158  
159 **Figure 2. Visualization of the three missing data simulation scenarios explored in this study.** On the  
160 left is Missing Completely at Random (MCAR) with a 40% missing rate. In the middle is Survey-Specific  
161 Missing Rate (SMR) with a 30% missing rate for Survey 1 and 50% missing rate for Survey 2. On the  
162 right is Blockwise Survey-Specific Missing Rate (BSMR) with a 30% missing rate for Survey 1 and 50%  
163 missing rate for Survey 2.

164  
165 ***Missing Completely at Random (MCAR)***

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

171  
172 ***Missing Not at Random (MNAR): Survey-Specific Missing Rate***

173 The second missing data simulation scenario is SMR, in which the proportion of missing values  
174 in each column is dependent on the survey type that it belongs to. SMR is tailored to mirror the  
175 missing rates in the full SPARK dataset by reusing the same proportions of missing values for each  
176 survey (**Table 1**).

177  
178 ***Missing Not at Random (MNAR): Blockwise Missingness with Survey-Specific Missing Rate***

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

185  
186 **3. Machine Learning Imputation**

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

192

### 193 *MICE*

194 This study used the MICE<sup>12</sup> (version 3.16.0) package in R which employs a multiple imputation  
195 model. It uses a concept called Fully Conditional Specification, in which each incomplete variable  
196 is imputed by a different model. It generates multiple imputed datasets that are averaged to retrieve  
197 the final imputed data. Since MICE employs a regression-based approach, hyperparameter tuning  
198 was not performed.

199

### 200 *KNN*

201 KNNImputer is a method in Python's Scikit-learn package<sup>22</sup> (version 0.22) and was used to study  
202 the KNN algorithm. KNNImputer predicts each sample's missing values by using the average  
203 value from the closest data points in the training set. Hyperparameter tuning was used to select the  
204 optimal value for the number of nearest neighbors used during imputation.

205

### 206 *MissForest*

207 MissForest<sup>15</sup> (version 1.5) is an R package which uses a Random Forest approach to impute  
208 missing values, building multiple decision trees to make predictions using the other remaining  
209 features. By averaging several classification or regression trees, MissForest employs out-of-bag  
210 error estimates and can capture complex, non-linear relationships. Hyperparameter tuning was  
211 used to select the optimal values for the number of trees and the maximum number of iterations.

212

### 213 *MIDAS*

214 MIDASpy<sup>23</sup> (version 1.3.1) is a Python package that was used to study the MIDAS algorithm. It  
215 introduces additional missing values into a given dataset and restores these values using an  
216 unsupervised neural network called a denoising autoencoder. Then, the resulting model is used to  
217 predict the values of the original missing data. Similar to MICE, MIDASpy generates multiple  
218 imputed datasets that are averaged to retrieve the final imputed data. Hyperparameter tuning was  
219 used to select the optimal values for the input drop, layer structure, and number of epochs.

220

## 221 **4. Evaluation of imputation performance**

222 For each missing data simulation scenario, we introduced missingness into the complete dataset  
223 10 different times as 10 separate trials. The values in **Table 1** correspond to the percentage of  
224 subject IDs in the full dataset (with missing values among participants with autism) who are not  
225 present in each specific survey. These missing rates were used when generating the missing  
226 datasets for the SMR and BSMR simulation scenarios.

227

228 The four models were used to impute the missing data, and these imputed values were compared  
229 with the true values in the preprocessed complete dataset. In each imputation trial, the RMSE  
230 values were calculated for each column using the `postResample` method from the `caret`  
231 package (Version 6.0-94) in R. The means of the RMSEs across all columns were aggregated to  
232 retrieve an overall RMSE. Then, these means were averaged across the 10 trials for each simulation

233 setting. This resulted in a mean overall RMSE for each simulation scenario. These error values  
234 were then compared for every simulation scenario between each imputation method.

235  
236 SCQ summary score, RBS-R summary score, and DCDQ summary score evaluate the social  
237 communication function, severity of repetitive behaviors, and motor functions respectively in  
238 study participants with autism. They were calculated based on corresponding questionnaires. The  
239 RMSE values of these specific mental and behavior summary scores were also compared between  
240 the four imputation methods across each simulation scenario.

241  
242 Lastly, the total computation time was assessed for the four imputation methods during the BSMR  
243 simulation scenario, which was chosen since it is closest in nature to missingness in real survey  
244 data.

245

## 246 **Results**

### 247 **1. Overview of full dataset and missingness patterns**

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



	Missing Rate			p-value
	Low missing rate (<20%)	Medium missing rate (20-80%)	High missing rate (>80%)	
<b>Number of Subjects</b>	37710 (32.2)	34067 (29.1)	45322 (38.7)	
<b>Sex (%)</b>				<0.001
Male	29460 (33.5)	24030 (27.3)	34412 (39.1)	
Female	8250 (28.3)	10037 (34.4)	10910 (37.4)	
<b>Age (%)</b>				<0.001
<2 years	456 (28.5)	636 (39.7)	509 (31.8)	
2-5 years	9773 (38.0)	6189 (24.1)	9726 (37.9)	
6-11 years	16511 (39.1)	9230 (21.9)	16463 (39.0)	
12-18 years	10966 (38.4)	6217 (21.7)	11401 (39.9)	
>18 years	4 (~0.0)	11795 (62.0)	7223 (38.0)	
<b>Race (%)</b>				<0.001
White	28727 (47.3)	17968 (30.0)	14093 (23.2)	
African American	2063 (37.8)	1373 (25.2)	2021 (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	4155 (48.3)	2203 (25.6)	2249 (26.1)	
Other	1654 (4.2)	11708 (30.0)	25767 (65.9)	

**Table 2. Sample characteristics by low (<20%), medium (20%-80%), and high (>80%) missing rate in SPARK.** Proportion of missing variables for each subject was calculated in the full dataset of this study containing 117,099 total ASD participants. Organized by different demographics including sex, age, and race.

254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271

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, only around 28% of female participants have low missing rates within this range.

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 towards extreme values, since around 39% have either low or high missing rates and 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 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 less than 25% of the

272 participants who were self-identified as White or “Multiple Races” reported high missing rates.  
273 Those who self-reported themselves as an “Other” race exhibit large amounts of missingness, since  
274 around 66% have missing rates larger than 80%.

275  
276 **2. Sample Characteristics of Complete Dataset and Simulation of Three Missingness**  
277 **Patterns**

278 To assess the imputation performance of the four popular missing data imputation methods (MICE,  
279 KNN, MissForest and MIDAS), we first obtained a preprocessed complete dataset with 15,196  
280 participants with autism (**Table 3**, details in Methods). Around 78% of participants with complete  
281 data are male and 22% are female. The male to female ratio is 3.5:1, which aligns with the sex  
282 ratio among subjects with autism in the general population. About half of the individuals with  
283 complete data are between 6-11 years of age. Only 0.4% of subjects are under 2 years of age while  
284 none are above 18. 79% of participants were self-identified as White. The category with the second  
285 largest number of participants is “Multiple Races” (10.9%), followed by African American (4.3%),  
286 followed by “Other” (3.5%), followed by Asian (2.2%). The number of participants who are Native  
287 American or Native Hawaiian are below 1%. In the preprocessed complete dataset, the SCQ, RBS-  
288 R, and DCDQ scores have average values of 21.72, 35.16, and 37.87 respectively.

289  
290

	<b>Number of Observations (Percentage) or Mean (Standard Deviation)</b>
<b>Number of Subjects</b>	15196
<b>Sex (%)</b>	
Male	11901 (78.3)
Female	3295 (21.7)
<b>Age (%)</b>	
<2 years	61 ( 0.4)
2-5 years	3029 (19.9)
6-11 years	8442 (55.6)
12-18 years	3664 (24.1)
>18 years	0 (0.0)
<b>Race (%)</b>	
White	11938 (78.6)
African American	656 ( 4.3)
Asian	331 ( 2.2)
Native American	71 ( 0.5)
Native Hawaiian	22 ( 0.1)
Multiple Races	1649 (10.9)
Other	529 ( 3.5)
<b>Summary Scores [mean (SD)]</b>	
SCQ Score	21.72 (7.09)
RBS-R Score	35.16 (20.50)
DCDQ Score	37.87 (12.73)

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

This table includes the number of observations and percentage breakdowns of sex, age, and race as well as means and standard deviations of the Social Communication Questionnaire (SCQ), Repetitive Behavior Scale-Revised (RBS-R), and Developmental Coordination Disorder Questionnaire (DCDQ) summary scores.

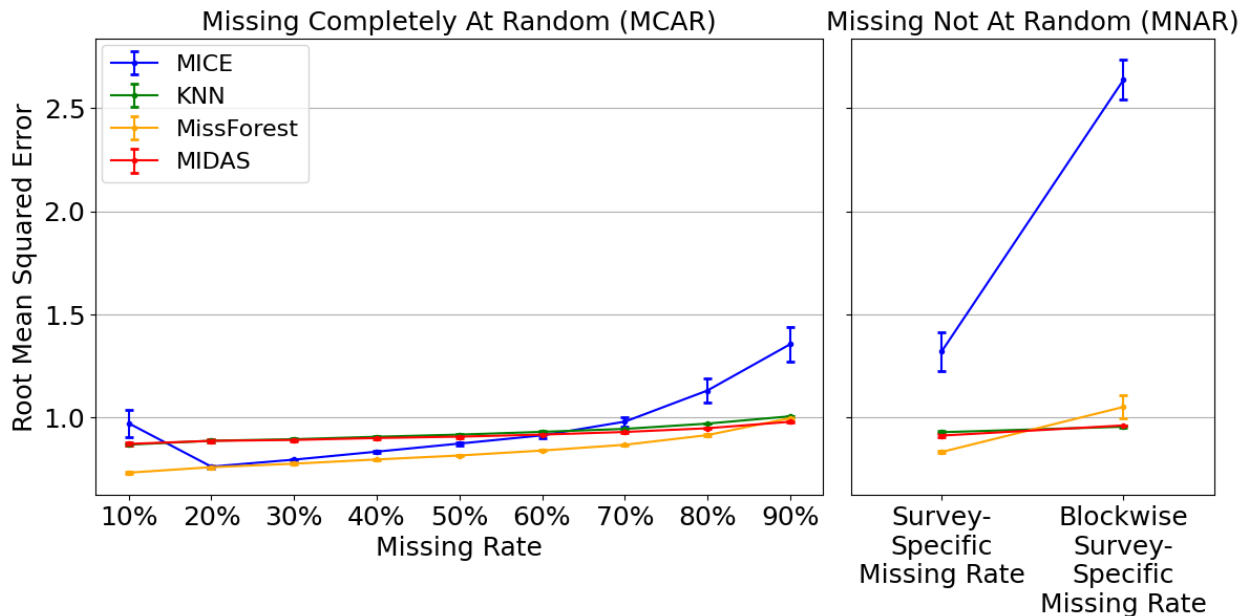
291  
292 To assess the performance of the missing data imputation methods, missing values were introduced  
293 to the preprocessed complete dataset with 15,196 participants with autism. First, to simulate the  
294 scenario on MCAR, a random subset of values across the entire dataset were converted to missing  
295 values. 10 incomplete datasets were generated for each missingness percentage (10%-90%).  
296 Second, to examine the performance of the imputation methods on MNAR patterns, 10 incomplete  
297 datasets were randomly generated for the SMR and BSMR simulation scenarios separately. When  
298 doing so, the missing rates in the original SPARK dataset were used (**Table 1**) to reflect the  
299 missingness distribution present in the real data.

300  
301  
302  
303  
304  
305

## 2. 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. Lower RMSE values correspond to higher accuracy in missing value imputation.

Evaluation of Imputation Performance based on Overall RMSE



306  
307  
308  
309  
310  
311

**Figure 3. Evaluation of imputation performance based on overall RMSE.** Values across the 10 trials using the MCAR simulation scenario (left). Overall RMSE values across the 10 MNAR trials in the Survey-Specific Missing Rate (SMR) and Blockwise Missingness with Survey-Specific Missing Rate (BSMR) simulation scenarios (right).

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 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% to 60% of random missingness but performed considerably worse than all other models for the remaining missing rates.

319

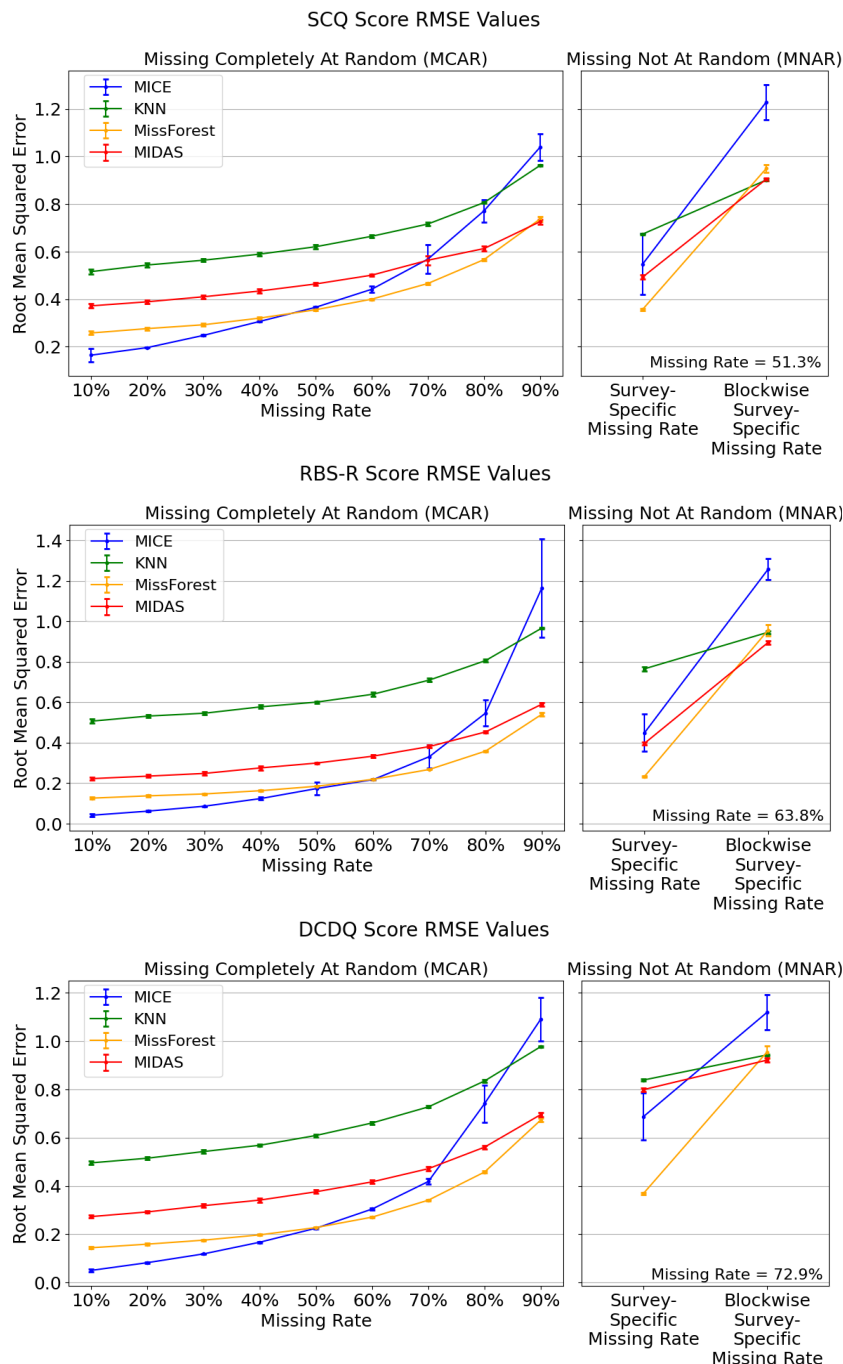
In the MNAR scenarios, all models exhibited an increase in imputation error in the BSMR scenario when compared to SMR. MissForest produced the lowest error rate in the SMR scenario, with an RMSE of 0.83, but did not perform as well during the BSMR scenario that simulated blockwise missingness. MissForest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

327

328 BSMR scenarios. Especially in the BSMR scenario, the RMSE value was significantly higher at  
329 2.64 with a relatively large standard deviation of 0.098.

330  
331 For every simulation scenario, the difference in imputation performance on overall RMSE between  
332 KNN and MIDAS was marginal. Both models produced very similar results throughout the  
333 experiment and for each simulation scenario besides BSMR, they typically performed slightly  
334 worse than MissForest.

335  
336 **3. Performance of Imputation on Mental and Behavioral Summary Scores**  
337 For every simulation scenario, the mean and standard deviations of RMSE values for the SCQ,  
338 RBS-R, and DCDQ scores were computed across the ten trials as displayed in **Figure 4**. The  
339 relative performance of the four models was generally consistent across the three summary scores.



340  
 341 **Figure 4. Imputation performance on summary scores from mental health surveys.** Root Mean  
 342 Squared Error (RMSE) values for imputing the Social Communication Questionnaire (SCQ) score across  
 343 the MCAR and MNAR trials (top). RMSE values for the Repetitive Behavior Scale-Revised (RBS-R)  
 344 score across the MCAR and MNAR trials (middle). RMSE values for the Developmental Coordination  
 345 Disorder Questionnaire (DCDQ) score across the MCAR and MNAR trials (bottom).

346  
 347 In the MCAR scenario, MissForest consistently outperformed KNN and MIDAS when imputing  
 348 all three summary scores. The MICE model exhibited a steep incline in error as the missing rate  
 349 was incremented. It performed the best until the missing rate was increased to 50%, after which it  
 350 was surpassed by the remaining models. MICE is ideal for lower rates of random missingness but

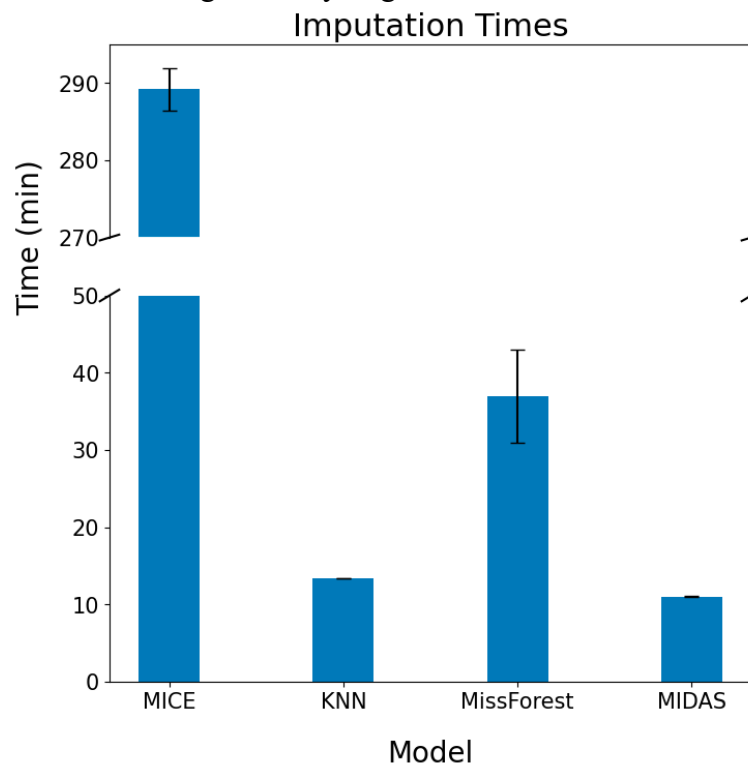
351 begins to perform exponentially worse as the rate gets larger. In fact, the MICE model produced  
352 the largest RMSE among the four methods at a 90% missing rate. For missing rates that are 50%  
353 and above, MissForest is the ideal model since it had the lowest errors among the four methods.  
354

355 The MissForest model performed the best in the SMR scenario. However, each method, especially  
356 MICE and MissForest, exhibited error rates that rose sharply when the missing values became  
357 blocked by survey type in the BSMR scenario. In the BSMR scenario, KNN and MIDAS exhibited  
358 the lowest error rates with MissForest performing slightly worse. MICE performed considerably  
359 worse than the remaining models in the BSMR scenario.  
360

#### 361 4. Computational Time

362 When comparing the computational times of the four models, the BSMR simulation scenario was  
363 used since this environment most closely resembles the missingness patterns in the real data when  
364 participants skip an entire survey in SPARK.  
365

366 As shown in **Figure 5**, MIDAS and KNN not only had similar overall error rates, but also exhibited  
367 comparable imputation times of around 10 to 13 minutes. MissForest had a median imputation  
368 time of slightly less than 30 minutes. On the other hand, MICE had a median imputation time of  
369 around 285 minutes, which was significantly larger than those of the remaining models.



370 **Figure 5.** Total imputation times (in minutes) and standard deviations of each model for the 10 trials in  
371 the Blockwise Missingness with Survey-Specific Missing Rate (BSMR) scenario. Total sample size is  
372 15,196.  
373  
374

#### 375 Discussion

376 The establishment of biobank databases has enabled the collection of self-reported mental and  
377 behavioral surveys at scale<sup>1-3</sup>. SPARK has gathered social and behavioral survey data from about  
378 100,000 individuals<sup>1</sup> and there is ongoing collection of more survey data on existing participants.  
379 UK Biobank has measurements on lifetime depressive disorder, cognitive function, attention, and  
380 impulsivity from about 150,000 participants<sup>2,24,25</sup>. All of Us also has strategic plans to collect  
381 mental and behavioral surveys at scale<sup>3</sup>. However, the data quality and statistical power are  
382 compromised by missing data. Recent advances in machine learning methods have inspired novel  
383 missing data imputation approaches with increased accuracy and computational efficiency<sup>12-15</sup>.  
384 Previous studies either have not reviewed these newly developed imputation methods or have not  
385 focused on assessing imputation accuracy in mental and behavioral surveys that exhibit blockwise  
386 missing structures<sup>17-21</sup>.

387  
388 Our study provided insights on the missingness pattern in SPARK, a large-scale cohort with  
389 autism, and assessed the imputation accuracy and computational time of four popular missing data  
390 imputation methods – MICE, KNN, MissForest and MIDAS. We did this by simulating three  
391 missingness scenarios in mental and behavioral surveys, including SCQ, RBS-R and DCDQ. We  
392 observed that 50%-70% of participants with autism did not complete SCQ, RBS-R and DCDQ  
393 surveys and the dataset exhibited blockwise missing structures. The missing rates also varied by  
394 sex, age, and race. Overall, KNN and MIDAS showed relatively stable performance with  
395 increasing missing rate in the MCAR scenario and slightly higher imputation error when blockwise  
396 missingness is introduced in the MNAR scenarios. The error rate increased more significantly in  
397 MICE and MissForest in both MCAR and MNAR scenarios, with a particularly notable surge in  
398 error rate for MICE when blockwise missing structures were introduced. When imputing SCQ,  
399 RBS-R and DCDQ summary scores in the MCAR scenario, MICE had the lowest error rate when  
400 the missing rate was low, while MissForest had the lowest error rate when the missing rate was  
401 high. However, in the presence of blockwise missingness in the MNAR scenario, MIDAS was  
402 consistently the best performing model across all three summary scores, with KNN and MissForest  
403 having similar or slightly higher error rates. Our results suggested that some models like MICE  
404 are sensitive to high missing rates and blockwise missing structures, while MIDAS and KNN may  
405 perform better in the overall dataset and specific summary scores in the presence of blockwise  
406 missingness. The average computational times for MIDAS and KNN to impute 15,196 subjects  
407 with blockwise missingness were about 10 minutes, about 35 minutes for MissForest, and about  
408 290 minutes for MICE. These results highlight the computational efficiency in machine learning  
409 imputation algorithms even in highly complex neural network models in MIDAS. Newly  
410 developed imputation models have better optimization in their algorithms and take advantage of  
411 parallel computing to reduce the computational time.

412  
413 Our results show the potential to impute missing data in large-scale databases with mental and  
414 behavioral surveys, especially imputing summary scores based on medical history and  
415 neurodevelopmental measures. When the data exhibits blockwise missingness, the imputation  
416 error increases but models such as MIDAS and KNN can still provide imputed results that are  
417 relatively stable and accurate. This shows that when a block of correlated variables in one survey  
418 is completely missing, other related surveys or medical history can also provide relevant  
419 information for imputation. The choice of imputation methods may depend on the overall missing  
420 rate and missingness patterns in a dataset.

421



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

429  
430 Missing data imputation is widely used in national surveys with mental and behavioral surveys.  
431 For example, the National Survey on Drug Use and Health (NSDUH) has been providing  
432 imputation-revised variables by the predictive mean neighborhood methods since 1999<sup>26</sup>. There is  
433 also the recent phenotype imputation model developed in the UK Biobank, which has shown  
434 increased power for genetic studies<sup>27</sup>. As biobanks and national surveys collect more large-scale  
435 data on mental and behavioral surveys, missing data imputation will produce more accurate  
436 imputed values and become an integral part of analysis to maximize the use of the data.

437  
438 Our study underscores the efficacy of advanced imputation techniques, such as MIDAS and KNN,  
439 in addressing missing data within large-scale mental and behavioral surveys. Our findings  
440 showcase that for similar databases with mental and behavioral surveys on autism, dementia and  
441 other disorders, machine learning-based imputation methods can be leveraged to effectively  
442 recover missing information. This study demonstrates that machine learning methods offer  
443 increased performance and faster computation times over traditional algorithms. The performance  
444 of these advanced imputation techniques demonstrates their potential to optimize analyses and  
445 advance research in mental and behavioral disorders.

446

## 447 **Figure Legends**

448 **Figure 1. Overview of workflow and study design.** a) The full dataset refers to the original  
449 data filtered to only include ASD participants. The preprocessed complete dataset refers to the  
450 original dataset after filtering to only include ASD participants, dropping incomplete rows,  
451 removing variables with extreme rates of missingness, and conducting one-hot-encoding on the  
452 categorical variables (which increases the number of variables). b) MCAR refers to the  
453 simulation scenario which randomly converts a specified fraction of the input dataset to missing.  
454 SMR refers to the simulation environment that is tailored to the missingness of the original  
455 dataset. BSMR refers to the simulation environment that is also tailored to the missingness of the  
456 original dataset, but converts all rows of a given column to missing at once. c) MICE is an  
457 imputation method that employs a series of regression models; MissForest is an imputation  
458 method that is based on random forests; MIDAS is an imputation method that uses denoising  
459 autoencoders; KNN is an imputation method that uses neighboring data points in the feature  
460 space. d) RMSE corresponds to Root Mean Squared Error.

461  
462 **Figure 2. Visualization of the three missing data simulation scenarios explored in this**  
463 **study.** On the left is Missing Completely at Random (MCAR) with a 40% missing rate. In the  
464 middle is Survey-Specific Missing Rate (SMR) with a 30% missing rate for Survey 1 and 50%  
465 missing rate for Survey 2. On the right is Blockwise Survey-Specific Missing Rate (BSMR) with  
466 a 30% missing rate for Survey 1 and 50% missing rate for Survey 2.

467

468 **Figure 3. Evaluation of imputation performance based on overall RMSE.** Values across the  
469 10 trials using the MCAR simulation scenario (left). Overall RMSE values across the 10 MNAR  
470 trials in the Survey-Specific Missing Rate (SMR) and Blockwise Missingness with Survey-  
471 Specific Missing Rate (BSMR) simulation scenarios (right).  
472

473 **Figure 4. Imputation performance on summary scores from mental health surveys.** Root  
474 Mean Squared Error (RMSE) values for imputing the Social Communication Questionnaire  
475 (SCQ) score across the MCAR and MNAR trials (top). RMSE values for the Repetitive  
476 Behavior Scale-Revised (RBS-R) score across the MCAR and MNAR trials (middle). RMSE  
477 values for the Developmental Coordination Disorder Questionnaire (DCDQ) score across the  
478 MCAR and MNAR trials (bottom).  
479

480 **Figure 5.** Total imputation times (in minutes) and standard deviations of each model for the 10  
481 trials in the Blockwise Missingness with Survey-Specific Missing Rate (BSMR) scenario. Total  
482 sample size is 15,196.  
483

## 484 **Table Legends**

485 **Table 1. Percentage of subjects who did not complete each individual survey among all**  
486 **117,099 participants with autism in SPARK.** Social Communication Questionnaire (SCQ),  
487 Repetitive Behavior Scale-Revised (RBS-R), and Developmental Coordination Disorder  
488 Questionnaire (DCDQ) are surveys commonly used to quantify the mental and behavioral  
489 functions at scale.  
490

491 **Table 2. Sample characteristics by low (<20%), medium (20%-80%), and high (>80%)**  
492 **missing rate in SPARK.** Proportion of missing variables for each subject was calculated in the  
493 full dataset of this study containing 117,099 total participants with autism. Organized by  
494 different demographics including sex, age, and race.  
495

496 **Table 3. Sample characteristics in the preprocessed complete dataset containing 15,196**  
497 **participants.** This table includes the number of observations and percentage breakdowns of sex,  
498 age, and race as well as means and standard deviations of the Social Communication  
499 Questionnaire (SCQ), Repetitive Behavior Scale-Revised (RBS-R), and Developmental  
500 Coordination Disorder Questionnaire (DCDQ) summary scores.

## 501 **Data availability**

502 SPARK Phenotype Dataset is accessible through application at SFARI Base  
503 (<https://base.sfari.org>)

## 504 **Code availability**

505 All software used in this study is publicly available. The code for simulations and analysis can be  
506 found at <https://github.com/AprilShuLab/MissingDataImputation>.

## 507 **Acknowledgements**

508 We are extremely grateful to the thousands of individuals and families who are participating in  
509 the SPARK. We thank the sites, staff and volunteers of the SPARK Clinical Site Network and  
510 SFARI for their invaluable contributions.

511 **Author contributions**

512 Preethi Prakash conducted the entire analysis and wrote the manuscript and Dr. Chang Shu  
513 supervised this work. Dr. Kelly Street, Dr. Shrikanth Narayanan and Dr. Yufeng Shen provided  
514 guidance on the methodology, and Dr. Bridget Fernandez offered insights on clinical relevance.

515

## 516 References

- 517 1 Feliciano, P. *et al.* SPARK: A US Cohort of 50,000 Families to Accelerate Autism  
518 Research. *Neuron* **97**, 488-493, doi:10.1016/j.neuron.2018.01.015 (2018).
- 519 2 Davis, K. A. S. *et al.* Mental health in UK Biobank - development, implementation and  
520 results from an online questionnaire completed by 157 366 participants: a reanalysis.  
521 *BJPsych Open* **6**, e18, doi:10.1192/bjo.2019.100 (2020).
- 522 3 Ramirez, A. H. *et al.* The All of Us Research Program: Data quality, utility, and diversity.  
523 *Patterns (N Y)* **3**, 100570, doi:10.1016/j.patter.2022.100570 (2022).
- 524 4 Chesnut, S. R., Wei, T., Barnard-Brak, L. & Richman, D. M. A meta-analysis of the social  
525 communication questionnaire: Screening for autism spectrum disorder. *Autism* **21**, 920-  
526 928, doi:10.1177/1362361316660065 (2017).
- 527 5 Hooker, J. L., Dow, D., Morgan, L., Schatschneider, C. & Wetherby, A. M. Psychometric  
528 analysis of the repetitive behavior scale-revised using confirmatory factor analysis in  
529 children with autism. *Autism Res* **12**, 1399-1410, doi:10.1002/aur.2159 (2019).
- 530 6 Van Damme, T., Vancampfort, D., Thoen, A., Sanchez, C. P. R. & Van Biesen, D.  
531 Evaluation of the Developmental Coordination Questionnaire (DCDQ) as a Screening  
532 Instrument for Co-occurring Motor Problems in Children with Autism Spectrum Disorder.  
533 *J Autism Dev Disord* **52**, 4079-4088, doi:10.1007/s10803-021-05285-1 (2022).
- 534 7 Jebb, A. T., Ng, V. & Tay, L. A Review of Key Likert Scale Development Advances:  
535 1995-2019. *Front Psychol* **12**, 637547, doi:10.3389/fpsyg.2021.637547 (2021).
- 536 8 Mack, C., Su, Z. & Westreich, D. in *Managing Missing Data in Patient Registries:  
537 Addendum to Registries for Evaluating Patient Outcomes: A User's Guide, Third Edition*  
538 (Agency for Healthcare Research and Quality (US), 2018).
- 539 9 Khan, S. I. & Hoque, A. S. M. L. SICE: an improved missing data imputation technique.  
540 *Journal of Big Data* **7**, 37, doi:10.1186/s40537-020-00313-w (2020).
- 541 10 Phiwhorm, K., Saikaew, C., Leung, C. K., Polpinit, P. & Saikaew, K. R. Adaptive multiple  
542 imputations of missing values using the class center. *Journal of Big Data* **9**, 52,  
543 doi:10.1186/s40537-022-00608-0 (2022).
- 544 11 de Goeij, M. C. M. *et al.* Multiple imputation: dealing with missing data. *Nephrology  
545 Dialysis Transplantation* **28**, 2415-2420, doi:10.1093/ndt/gft221 (2013).
- 546 12 van Buuren, S. & Groothuis-Oudshoorn, K. mice: Multivariate Imputation by Chained  
547 Equations in R. *Journal of Statistical Software* **45**, 1 - 67, doi:10.18637/jss.v045.i03  
548 (2011).
- 549 13 Azur, M. J., Stuart, E. A., Frangakis, C. & Leaf, P. J. Multiple imputation by chained  
550 equations: what is it and how does it work? *Int J Methods Psychiatr Res* **20**, 40-49,  
551 doi:10.1002/mpr.329 (2011).
- 552 14 Taunk, K., De, S., Verma, S. & Swetapadma, A. A Brief Review of Nearest Neighbor  
553 Algorithm for Learning and Classification. (2019).
- 554 15 Stekhoven, D. J. & Bühlmann, P. MissForest—non-parametric missing value imputation  
555 for mixed-type data. *Bioinformatics* **28**, 112-118, doi:10.1093/bioinformatics/btr597  
556 (2011).
- 557 16 Lall, R. & Robinson, T. The MIDAS Touch: Accurate and Scalable Missing-Data  
558 Imputation with Deep Learning. *Political Analysis* **30**, 179-196, doi:10.1017/pan.2020.49  
559 (2022).
- 560 17 Shrive, F. M., Stuart, H., Quan, H. & Ghali, W. A. Dealing with missing data in a multi-  
561 question depression scale: a comparison of imputation methods. *BMC Medical  
562 Research Methodology* **6**, 57, doi:10.1186/1471-2288-6-57 (2006).
- 563 18 Peyre, H., Leplège, A. & Coste, J. Missing data methods for dealing with missing items  
564 in quality of life questionnaires. A comparison by simulation of personal mean score, full  
565 information maximum likelihood, multiple imputation, and hot deck techniques applied to

- 566 the SF-36 in the French 2003 decennial health survey. *Quality of Life Research* **20**, 287-  
567 300, doi:10.1007/s11136-010-9740-3 (2011).
- 568 19 Emmanuel, T. *et al.* A survey on missing data in machine learning. *Journal of Big Data* **8**,  
569 140, doi:10.1186/s40537-021-00516-9 (2021).
- 570 20 Xu, X. *et al.* The ability of different imputation methods for missing values in mental  
571 measurement questionnaires. *BMC Medical Research Methodology* **20**, 42,  
572 doi:10.1186/s12874-020-00932-0 (2020).
- 573 21 Croy, C. D. & Novins, D. K. Methods for addressing missing data in psychiatric and  
574 developmental research. *J Am Acad Child Adolesc Psychiatry* **44**, 1230-1240,  
575 doi:10.1097/01.chi.0000181044.06337.6f (2005).
- 576 22 Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *the Journal of machine*  
577 *Learning research* **12**, 2825-2830 (2011).
- 578 23 Lall, R. & Robinson, T. Efficient Multiple Imputation for Diverse Data in Python and R:  
579 MIDASpy and rMIDAS. *Journal of Statistical Software* **107**, 1 - 38,  
580 doi:10.18637/jss.v107.i09 (2023).
- 581 24 Fawns-Ritchie, C. & Deary, I. J. Reliability and validity of the UK Biobank cognitive tests.  
582 *PLOS ONE* **15**, e0231627, doi:10.1371/journal.pone.0231627 (2020).
- 583 25 Schweren, L. J. S. *et al.* Diet, Physical Activity, and Disinhibition in Middle-Aged and  
584 Older Adults: A UK Biobank Study. *Nutrients* **13**, 1607 (2021).
- 585 26 Grau, E., Frechtel, P., Odom, D. & Painter, D. in *2004 Proceedings of the Section on*  
586 *Survey Research Methods*.
- 587 27 An, U. *et al.* Deep learning-based phenotype imputation on population-scale biobank  
588 data increases genetic discoveries. *Nature Genetics* **55**, 2269-2276,  
589 doi:10.1038/s41588-023-01558-w (2023).
- 590