
Tabular data imputation: quality over quantity

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Tabular data imputation algorithms allow to estimate missing values and use
2 incomplete numerical datasets. Current imputation methods minimize the error
3 between the unobserved ground truth and the imputed values. We show that this
4 strategy has major drawbacks in the presence of multimodal distributions, and
5 we propose to use a qualitative approach rather than the actual quantitative one.
6 We introduce the kNNxKDE algorithm: a hybrid method using chosen neighbors
7 (k NN) for conditional density estimation (KDE) tailored for data imputation. We
8 qualitatively and quantitatively show that our method preserves the original data
9 structure when performing imputation. This work advocates for a careful and
10 reasonable use of statistics and machine learning models by data practitioners.

11 1 Introduction

12 Big data is often referred to as the "gold of the 21st century". But with ubiquitous large databases,
13 missing data are a pervasive problem. They can introduce a bias, lead to wrong conclusions, or even
14 prevent from using data analysis tools that require complete datasets.

15 To mitigate this issue, data imputation algorithms have been developed. From the straightforward
16 mean/mode imputation to recent artificial neural networks (ANN) models, a wide range of tools
17 are available to impute incomplete datasets. This study focuses on tabular datasets, i.e. numerical
18 data arranged in rows and columns in a form of a matrix. For tabular datasets, recent benchmarks
19 argue that complex imputation methods do not perform better than simple traditional algorithms
20 [Bertsimas et al., 2018, Poulos and Valle, 2018, Jadhav et al., 2019, Woznica and Biecek, 2020, Jäger
21 et al., 2021]. In particular, the consensus is that the k NN-Imputer [Troyanskaya et al., 2001] and
22 MissForest [Stekhoven and Bühlmann, 2012], in spite of being traditional and simple algorithms,
23 generally perform better over a large range of datasets in various missing data scenarios.

24 Data may be missing because it was not recorded, the record has been lost, degraded, or the data may
25 also be censored. Missing data scenarios are usually classified into three types [Little and Rubin,
26 2014]: missing completely at random (MCAR), missing at random (MAR) and missing not at random
27 (MNAR). In MCAR the missing data mechanism is assumed independent of the dataset. In MAR,
28 the missing data mechanism is assumed to only dependent on the observed variables. The MNAR
29 scenario encompasses all other possible scenarios: the reason why data is missing may depend on the
30 missing value itself. Most comparisons focus on the MCAR scenario.

31 Tabular data imputation methods have always been evaluated using the RMSE between the estimated
32 value and the ground truth. The higher the mean RMSE, the poorest the imputation method. This
33 approach is of course intuitive, but is too restrictive for multimodal datasets: it assumes that for a set
34 of observed variables, there exists only a unique answer to recover. For multimodal datasets, density
35 estimation methods like the familiar Kernel Density Estimation (KDE) [Rosenblatt, 1956, Parzen,
36 1962], appear of interest for data imputation. But despite some attempts [Titterington and Mill, 1983,

37 Leibrandt and Günnemann, 2018], density estimation methods do not handle well observations with
38 missing values.

39 In this paper, we propose to step back and look at simple datasets to demonstrate that current
40 approaches for data imputation have serious shortcomings. To tackle them, we introduce a local
41 density estimator tailored for data imputation. By leveraging the convenient properties of the k NN-
42 Imputer and KDE, we develop k NNxKDE: a simple yet efficient algorithm for stochastic local data
43 imputation. We visually show that our method performs better than standard methods, and evaluate
44 the performances using the likelihood when available. We provide the code and the data used in
45 this work for reproducibility. Interested readers may experiment with the hyperparameters of our
46 algorithm.

47 2 Current methods perform poorly for multimodal dataset

48 This section demonstrates that conventional data imputation methods provide poor imputation with
49 basic multimodal datasets. For this purpose, we generate three simple two-dimensional datasets and
50 visually assess the imputation performances of four standard methods.

51 2.1 Three simple datasets

52 The first dataset is a bijection. x_1 is sampled from a mollified uniform distribution on $[0, 1]$ with
53 standard deviation $\sigma = 0.05$. Then $x_2 = x_1 + \varepsilon$, where $\varepsilon \sim N(0, 0.1)$.

54 The second dataset is a surjection, using a sine wave: $x_1 = 4\pi u$, where u is sampled from a mollified
55 distribution on $[0, 1]$ with standard deviation $\sigma = 0.05$. Then $x_2 = \sin x_1 + \varepsilon$, where $\varepsilon \sim N(0, 0.2)$.
56 The surjection allows to show that most imputation algorithms perform well in the unambiguous case
57 (when x_2 is missing), but not with multimodal distributions (when x_1 is missing).

58 Finally, Dataset 3 displays a ring. It has been generated in polar coordinates: $\theta \sim \mathcal{U}[0, 2\pi]$ and
59 $r = 1.0 + \varepsilon$, where $\varepsilon \sim N(0, 0.1)$. Euclidean coordinates are $x_1 = r \cos \theta$ and $x_2 = r \sin \theta$.

60 All three datasets have $N = 500$ observations and are plotted in Figure 1. The code used for
61 generation and the datasets themselves are provided in supplementary materials. We have used a
62 mollified uniform distribution for x_1 in Datasets 1 and 2 to prevent from zero likelihood computation
63 problems at the edges of the uniform distribution.

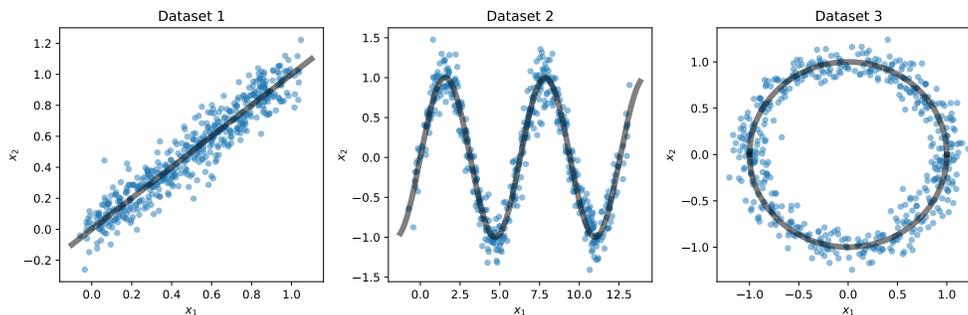


Figure 1: Three basic synthetic datasets with $N = 500$ observations. Dataset 1 is a bijection, Dataset 2 is a surjection, and Dataset 3 uses polar coordinates (not a function in the euclidean space).

64 2.2 Four standard data imputation methods

65 Here, we present the four data imputation methods used in this work: the k NN-Imputer, MissForest,
66 MICE and GAIN. This choice is of course arbitrary, but illustrates well the current state of affairs
67 regarding tabular data imputation [Bertsimas et al., 2018, Poulos and Valle, 2018, Yoon et al., 2018,
68 Jadhav et al., 2019, Woznica and Biecek, 2020, Jäger et al., 2021]

- 69 • The k NN-Imputer [Troyanskaya et al., 2001] computes distances between pairs of obser-
70 vations using a Euclidean distance that can handle missing values (called nan-Euclidean

71 distance). It imputes missing values by looking at one column at a time and averaging over
72 the k nearest neighbors that have an observed value for that column. Therefore, different
73 neighbors can be used to impute two missing entries in the same observation. One needs to
74 tune the hyperparameter k for the number of neighbors. The scientific consensus puts the
75 k NN-Imputer often on par with MissForest as for the best tabular data imputation method.

- 76 • MissForest [Stekhoven and Bühlmann, 2012] is an iterative imputation algorithm. It begins
77 by filling all missing values with initial estimates (e.g. the column mean), and then loops
78 through all columns, one at a time, performing a regression of that specific column onto all
79 other columns using Random Forests. It stops when the imputed dataset is stable enough
80 (following a user-defined threshold). The number of trees has to be tuned. MissForest has
81 shown great flexibility and successful data imputation results.
- 82 • MICE stands for Multiple Imputation Chained Equations [van Buuren and Groothuis-
83 Oudshoorn, 2011]. Similar to MissForest, it is an iterative imputation algorithm that uses
84 a regressor (linear regressions for MICE) to predict each column successively after filling
85 all missing entries with initial guesses. This algorithm has no hyperparameter to optimize.
86 MICE has shown good imputation results and is appreciated for its simplicity and absence
87 of hyperparameter tuning, but it fails at capturing non-linear dependencies.
- 88 • Finally, GAIN is a GAN neural network tailored for tabular data imputation which claims
89 state-of-the-art imputation results [Yoon et al., 2018]. GAIN smartly revisits the GAN
90 architecture by working with individual cells rather than whole observations. It has benefited
91 from a lot of attention for tabular data imputation. However, recent benchmarks show
92 that its performances are mediocre in practice [Jäger et al., 2021]. GAIN has several
93 hyperparameters to tune: batch size, hint rate (amount of correct labels provided to the
94 discriminator), number of training iterations, and weight parameter α for the generator loss
95 (balances RMSE loss for the observed cells and adversarial loss for the generated cells). We
96 decide to follow the authors’ recommendations and fix: batch size $N_{\text{batch}} = 128$, hint rate
97 $r_h = 0.9$ and $\alpha = 100$. We only optimize the number of iterations.

98 2.3 Imputation results

99 We introduce missing values for each dataset in a MCAR scenario with 20% missing rate. If an
100 observation has both features removed, we repeat the process until at least one feature is present.
101 After missing values have been injected, we normalize the dataset in the range $[0, 1]$ using the minimum
102 and maximum value of each feature.

For each data imputation algorithm and for each dataset, we perform a grid search of the hyperparam-
eter that best minimizes the normalized RMSE (NRMSE):

$$\text{NRMSE} = \sqrt{\frac{1}{N_{\text{miss}}} \sum_{i=1}^n \sum_{j=1}^d (x_{ij} - \hat{x}_{ij})^2 m_{ij}}$$

103 where $m_{ij} = 1$ if cell (i, j) is missing ($m_{ij} = 0$ otherwise) and $N_{\text{miss}} = \sum_{i=1}^n \sum_{j=1}^d m_{ij}$ is the
104 total number of missing entries in the dataset. The best hyperparameters, presented in Table 1, are
105 used to impute each dataset one more time. The optimized imputation results are plotted in Figure 2.

Table 1: Hyperparameter search results for each imputation method and dataset

	Data imputation method			
	k NN-Imputer	MissForest	MICE	GAIN
Dataset 1	$k = 30$ neighbors	$N_{\text{trees}} = 10$	X	$N_{\text{iter}} = 500$
Dataset 2	$k = 30$ neighbors	$N_{\text{trees}} = 30$	X	$N_{\text{iter}} = 200$
Dataset 3	$k = 75$ neighbors	$N_{\text{trees}} = 30$	X	$N_{\text{iter}} = 100$

106 We believe that Figure 2 provides meaningful insight regarding the current state of tabular data
107 imputation. The scientific consensus is that the k NN-Imputer and MissForest provide overall better
108 data imputation quality, which is somewhat recovered here. MICE uses linear regression between
109 features and cannot capture non-linear dependencies. Despite its flexible architecture, GAIN do not
110 recover missing values, even for Dataset 1. GAIN is hard to train properly.

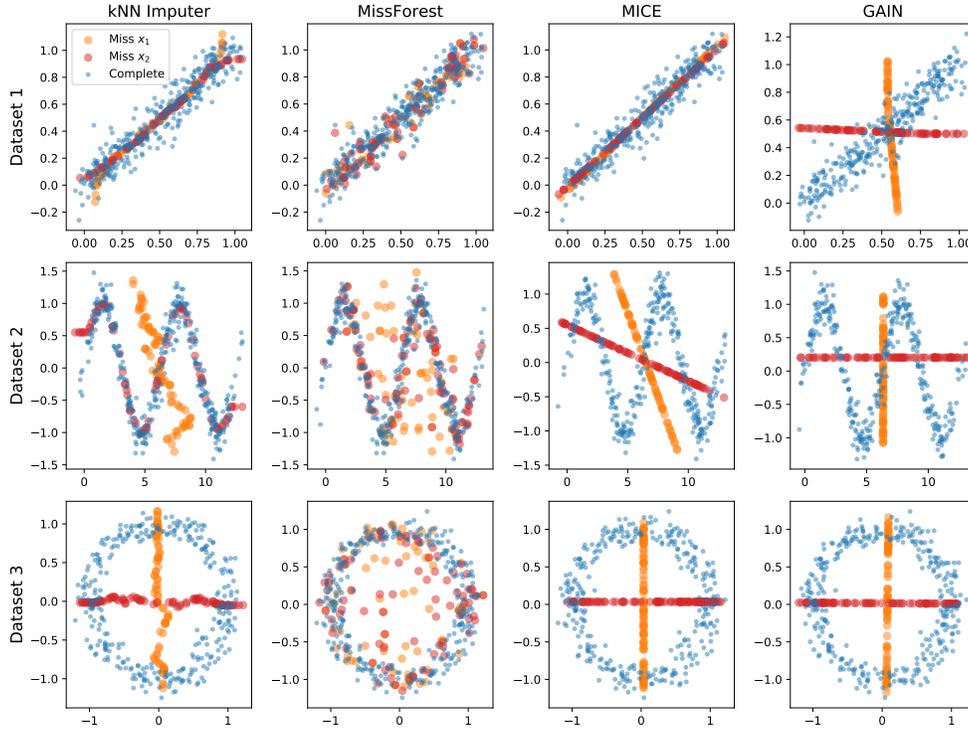


Figure 2: Imputation results for the three synthetic datasets by the four selected imputation methods with optimized hyperparameters. Blue dots correspond to complete observations, orange dots have observed x_2 but imputed x_1 , and red dots have observed x_1 but imputed x_2 . The k NN-Imputer, MissForest and MICE perform well on Dataset 1. The k NN-Imputer and MissForest can impute x_2 for Dataset 2, but they cannot impute x_1 . No method can properly impute Dataset 3. GAIN provides the worst imputation results and cannot even impute Dataset 1.

111 Both the k NN-Imputer and MissForest average over several predictions. This is why the imputation
 112 of x_1 in Dataset 2 lies between the two sine waves, and imputations for both x_1 and x_2 in Dataset
 113 3 are inside the ring. While averaging over several predictions often lead to better estimates, this
 114 strategy deteriorates the imputation quality if the missing value distribution is not unimodal.

115 MICE performs imputation by assuming linear dependency between features in the dataset. It is
 116 therefore no surprise if MICE can very well impute Dataset 1 but fails at imputing Dataset 2 and
 117 Dataset 3. Once the MICE algorithm has converged, the imputed orange and red dots follow almost
 118 perfectly the center of mass of all points in the dataset.

119 GAIN provides surprisingly disappointing imputation results. While ANNs are flexible models, the
 120 generator and the discriminator of GAIN fail to capture the non-linear relationship between x_1 and
 121 x_2 in all three datasets. Because of its innovative and complex framework, GAIN suffers from a
 122 complicated training process, which leads to bad imputation results. We have tried to train GAIN
 123 several times with various hyperparameters, but always end up with similar imputation quality.

124 3 kNNxKDE

125 To address the issues presented in Section 2, we propose a local stochastic imputer using kernel
 126 density estimation with Gaussian kernels. We adapt the KDE algorithm to missing data settings: only
 127 the conditional density of missing features given the observed features is estimated.

128 We use a methodology analogous to the k NN-Imputer to look for neighbors, but we work with
 129 missing patterns instead of working column by column. The reason of this choice is that working
 130 with one column at a time may lead to incoherent imputations as the selected neighbors for different

131 columns are different. Therefore, some imputed observations may be incompatible with the dataset
 132 structure. For a dataset with D columns, we have up to $2^D - 2$ possible missing patterns. Indeed,
 133 each cell may either be missing or not (hence 2^D choices) but we do not account for complete cases
 134 (nothing to impute) and completely unobserved cases (without even an observed cell).

For each pair of observations in the normalized dataset, we compute the distance d_{ij} using the
 nan-Euclidean distance [Dixon, 1979]:

$$d_{ij} = \sqrt{\frac{D}{|\mathcal{D}_{\text{obs}}|} \sum_{k \in \mathcal{D}_{\text{obs}}} (x_{ik} - x_{jk})^2}$$

where D is the total number of columns in the dataset, $\mathcal{D}_{\text{obs}} = \{k \in \llbracket 1, D \rrbracket \mid m_{ik} = m_{jk} = 1\}$ is
 the set of indices for commonly observed features in observations i and j and $|\mathcal{D}_{\text{obs}}|$ is its cardinality.
 These pairwise distances are then passed to a softmax function in order to define probabilities:

$$p_{ij} = \frac{e^{-\tau d_{ij}}}{\sum_j e^{-\tau d_{ij}}}$$

135 We use the "soft" version of the k NN algorithm, and introduce the temperature hyperparameter τ .
 136 Instead of selecting a fixed number of neighbors per observation, we use a neighborhood where
 137 nearest neighbors have stronger weights. In a similar fashion as Frosst et al. [2019], the notion of
 138 temperature controls the tightness of each observation's neighborhood.

139 Given a missing pattern, we first select all data to impute and potential donors. Data to impute is
 140 the subset of data which has the current missing pattern, and potential donors are the subset of data
 141 where at least all columns in the current missing pattern are observed. For an incomplete observation
 142 i in the subset of data to impute, p_{ij} is the probability of choosing observation j from the subset of
 143 potential donors. We have $\sum_j p_{ij} = 1$. Algorithm 1 shows the pseudo-code of the kNNxKDE.

144 The kNNxKDE has three hyper-
 parameters. The temperature τ
 for the softmax probabilities, the
 (shared) standard deviation h of
 the Gaussian kernels, and the num-
 ber N_{draws} of total sampled neigh-
 bors. The temperature τ controls
 the breadth of the selected neigh-
 borhood. The standard deviation
 h corresponds to the width of the
 Gaussian kernels. The effects of τ
 and h are discussed in Section 4.
 The last hyperparameter is the num-
 ber N_{draws} of imputation samples
 to be returned. It determines the
 resolution of the estimated density.
 Besides the obvious computational
 resources, there are no drawbacks
 to setting a high number of imputa-
 tion samples N_{draws} .

Algorithm 1: Pseudo-code for the kNNxKDE

Data: The incomplete dataset X
 min/max normalization;
for each missing pattern do
 $X_{\text{imp}} \leftarrow \text{data_to_impute};$
 $X_{\text{don}} \leftarrow \text{potential_donors};$
 $d_{ij} \leftarrow \text{nanEuclidDist}(X_{\text{imp}}, X_{\text{don}});$
 if d_{ij} *is NaN* **then**
 $d_{ij} \leftarrow \infty;$
 end
 $p_{ij} \leftarrow \text{softmax}(-\tau d_{ij});$
 for each row in X_{imp} **do**
 $r \leftarrow \text{sample } N_{\text{draws}} \text{ indices in } X_{\text{don}} \text{ with prob } p_{ij};$
 $e \leftarrow \text{sample } N_{\text{draws}} \text{ from } e \sim \mathcal{N}(0, h);$
 $\text{imputation_samples} \leftarrow X_{\text{don}}[r] + e;$
 end
end
 min/max renormalization;
Return: $\text{imputations_samples}$

145 **4 Results on synthetic datasets**

146 In Subsection 4.1, we show the performances of the kNNxKDE on the three artificial datasets and we
 147 discuss the effect of the hyperparameters τ and h . In Subsection 4.2, we use the log-likelihood of the
 148 imputed sample as an attempt to quantify imputation quality. We show that, for multimodal datasets,
 149 using the likelihood is more appropriate than the RMSE. All experiments use the MCAR setting to
 150 artificially introduce missing data with 20% missing rate.

151 **4.1 Qualitative evaluation of the kNNxKDE algorithm**

152 We show that the proposed method provides imputation samples that preserve the structure of the
 153 original dataset. For now, we fix the hyperparameters of the kNNxKDE at their default values:
 154 $h = 0.03$, $\tau = 50.0$ and $N_{\text{draws}} = 10000$. Figure 3 shows the imputation with a sub-sampling
 155 size $N_{\text{ss}} = 10$. The sub-sampling size is only used to show the variability in the imputation results
 156 by sampling several times. If x_1 is missing, we sample N_{ss} possible values given x_2 (the orange
 157 horizontal trails of dots), and if x_2 is missing, we draw N_{ss} possible estimates given x_1 (the red
 158 vertical trails of dots).

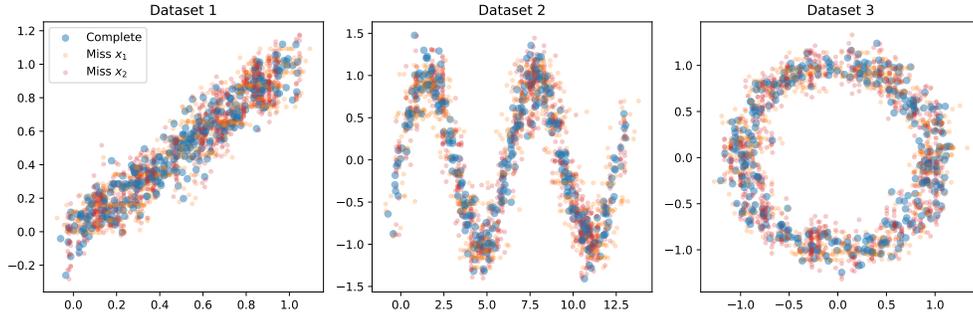


Figure 3: Several imputation results from the kNNxKDE algorithm. Each missing entry has been imputed $N_{\text{ss}} = 10$ times to show the variability of the estimates. The imputed values match with the structure of the observed data (larger blue dots).

159 Another way to visualize the distribution of the conditional distribution for each missing value
 160 is to look at the univariate density provided by the kNNxKDE algorithm. For each dataset, we
 161 have selected two observations: one with missing x_1 and one with missing x_2 . Figure 4 shows six
 162 univariate densities returned by the kNNxKDE algorithm with default hyperparameters values. In the
 163 upper left corner of each panel, the observed value is shown for reference. On each panel, a thick
 164 dashed line indicates the (unknown) ground truth. We see that the ground truth always falls in one of
 165 the modes of the estimated imputation density.

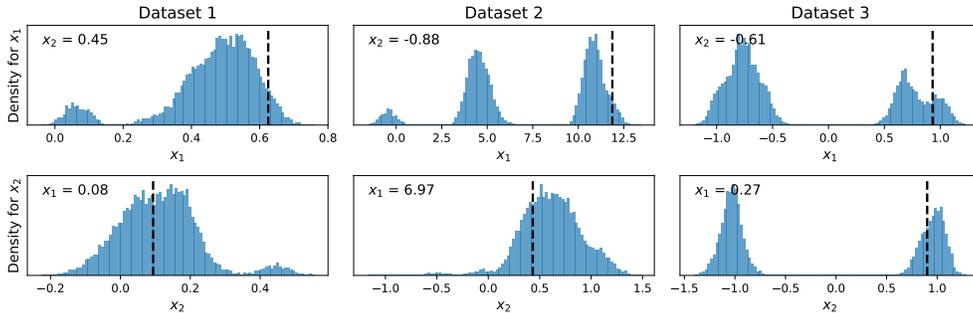


Figure 4: Example of conditional density distributions from the kNNxKDE algorithm with default hyperparameter values. Each histogram has $N_{\text{draws}} = 10000$ samples. Thick dashed lines correspond to the (unobserved) ground truth and the observed value is in the upper-left corner.

166 For Dataset 2, when x_1 is missing (upper middle panel of Figure 4), the kNNxKDE returns a
 167 multimodal distribution. Indeed, given the observed $x_2 = -0.88$, three separate ranges of values
 168 could correspond to the missing x_1 . Similarly, Dataset 3 shows bimodal distributions both for x_1 or
 169 x_2 , corresponding to the two possible ranges of values allowed by the ring structure.

170 We now focus on Dataset 2 to experiment with the hyperparameters h and τ . Figure 5 shows how the
 171 imputation quality changes when we vary the softmax temperature τ , and the effects of the Gaussian
 172 kernel bandwidth h are shown in Figure 6.

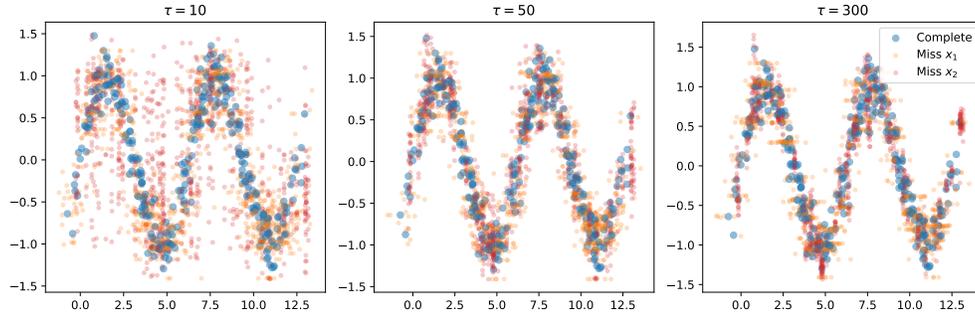


Figure 5: Evolution of the imputation quality as the softmax temperature τ varies. The Gaussian kernel bandwidth is fixed at $h=0.03$. We see that if τ is too low, the imputation has a large variance. If τ is too high, the imputation could be biased.

173 The value of the softmax temperature τ plays an important role in the data imputation quality, as can
 174 be seen in Figure 5. Recall that τ constrains the neighborhood range for each observation. The lower
 175 τ , the looser the neighborhood, and irrelevant observations could be sampled. This results in a large
 176 scatter (leftmost panel). Conversely, the higher τ , the tighter the neighborhood. Missing values will
 177 be imputed using very few other observations and multimodality can be overlooked. This can be seen
 178 on the rightmost panel, where the sampling variability is only due to the Gaussian kernel bandwidth.
 179 Tuning τ means finding a good balance in the bias/variance tradeoff.

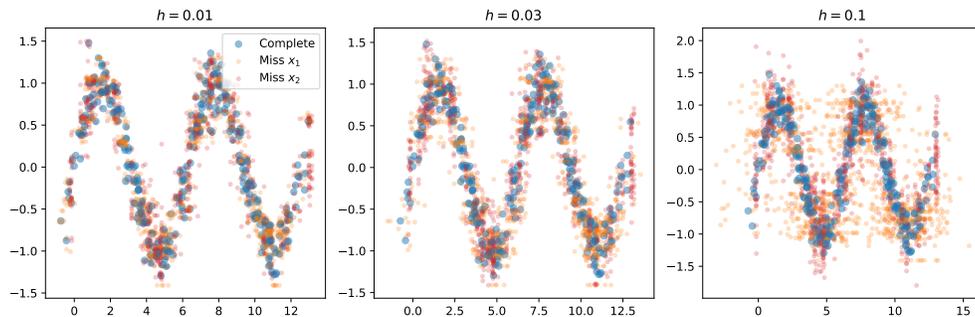


Figure 6: Change in the imputation quality when the Gaussian kernel bandwidth h varies. The softmax temperature is fixed at $\tau = 50$. We see that if h is too low, the imputation sample is very close to the observed data. If h is too high, the imputation sample is too scattered.

180 Now, the kernel bandwidth h controls the amount of fit to the observed data (c.f. Figure 6). The
 181 lower h , and the closer to the observed data the imputation sample will be. This can result in spiky
 182 univariate distributions. In the limit where $h = 0.0$, the conditional distribution for each missing value
 183 becomes a multinomial distribution with probability given by the softmax function computed with
 184 the pairwise distances. On the contrary, the higher h and the higher the variability of the imputation
 185 sample. Unlike τ , a bandwidth h too narrow does not mean that multimodality will be overlooked.
 186 With low h , the univariate distribution for a multimodal conditional probability will show distinct
 187 pronounced peaks. If h is too high, the different modes may collapse into a larger distribution with
 188 high variance.

189 4.2 The log-likelihood to measure imputation quality

190 Here, we compute the normalized RMSE (NRMSE) for the three datasets after imputation with all
 191 standard methods and the kNNxKDE algorithm. We compare the NRMSE with the log-likelihood
 192 score, which we can also compute since we know the generative process of the synthetic datasets.
 193 When performing a single imputation with the kNNxKDE algorithm, we draw a unique random
 194 sample from the resulting imputation distribution.

195 For each dataset and each imputation method, we repeat 100 times the following process: we introduce
 196 missing values, normalize the dataset, impute with the selected method using best hyperparameters
 197 (c.f. Table 1) and compute the NRMSE. Table 2 shows the mean and the standard deviation of the
 198 NRMSE. As already discussed in Section 2, the k NN-Imputer, MissForest and MICE have a low
 199 RMSE for Dataset 1, meaning that these methods recover well missing values. Larger NRMSEs for
 200 Datasets 2 and 3 quantify the poorer imputation quality. GAIN has a large RMSE, even for Dataset 1,
 201 as it could be anticipated from Section 2.

Table 2: Normalized RMSE for the three datasets with all imputation methods. kNNxKDE does not perform particularly well in terms of minimizing the NRMSE.

	Data imputation method				
	k NN-Imputer	MissForest	MICE	GAIN	kNNxKDE
Dataset 1	0.075 \pm 0.005	0.096 \pm 0.005	0.075 \pm 0.004	0.228 \pm 0.026	0.111 \pm 0.006
Dataset 2	0.192 \pm 0.011	0.252 \pm 0.019	0.250 \pm 0.009	0.271 \pm 0.023	0.267 \pm 0.017
Dataset 3	0.295 \pm 0.010	0.374 \pm 0.022	0.294 \pm 0.010	0.309 \pm 0.027	0.419 \pm 0.024

202 The kNNxKDE does not perform well with the RMSE. It has the largest NRMSEs, if we disregard
 203 GAIN. The justification we provide is that the kNNxKDE is not designed to accurately recover
 204 missing values. When performing a single imputation, the kNNxKDE algorithm selects a unique
 205 sample from the resulting imputation distribution. This is equivalent to selecting a single neighbor
 206 with the softmax probabilities – which may not even be the closest neighbor – and using a noisy copy
 207 of its observed values for imputation. This is an audacious choice, while the other imputation methods
 208 look for an optimal compromise. For multimodal distributions, sampling with the kNNxKDE cannot
 209 guarantee that we sample from the mode where the ground truth lies. For Dataset 3, where kNNxKDE
 210 shows the highest NRMSE, the imputation may be completely off (i.e., on the other side of the ring).

211 We now compute the log-likelihood of the resulting imputed sample. Like with the NRMSE, for
 212 each dataset and each imputation method, we repeat 100 independent experiments with the best
 213 hyperparameters. The imputed data are renormalized back to their original range to compute the
 214 log-likelihood of the imputed samples. Table 3 shows the mean and the standard deviation of the
 215 log-likelihood.

Table 3: Mean and standard deviation of the log-likelihood for the three datasets with all imputation methods. The first column shows the log-likelihood of the original sample for reference.

	Ref.	Data imputation method				
		k NN-Imputer	MissForest	MICE	GAIN	kNNxKDE
Dataset 1	425	494 \pm 9	450 \pm 14	495 \pm 11	-234 \pm 231	408 \pm 15
Dataset 2	79	-2214 \pm 299	-525 \pm 150	-2691 \pm 261	-1482 \pm 600	-54 \pm 33
Dataset 3	-481	-2251 \pm 196	-893 \pm 117	-2361 \pm 209	-2117 \pm 319	-509 \pm 15

216 This time, kNNxKDE performs best for Datasets 2 and 3. For Dataset 1, the k NN-Imputer, MissForest
 217 and MICE have a larger log-likelihood than the original sample because these methods average over
 218 several predictions and therefore remove the variability in their predictions: the imputed sample is
 219 very close to the ground truth and shows a high likelihood under the generative model (c.f. Figure 2).
 220 The log-likelihood of the imputed samples by GAIN is poor regardless of the dataset. MissForest
 221 shows interestingly decent results as it benefits from the iterative imputation mechanism and the
 222 random forest flexibility to capture non-linear dependency (unlike MICE).

223 With the log-likelihood as the new evaluation metric, the kNNxKDE now provides the best imputed
 224 samples. Each imputed observation may be far from its ground truth – hence the large NRMSE in
 225 Table 2, but it conforms to the data structure – hence the large log-likelihood in Table 3.

226 5 Discussion

227 We have shown the limits of the RMSE for data imputation problems, and have introduced a new
 228 data imputation method. In this last section, we talk about the limitations and the strengths of the

229 kNNxKDE algorithm, and summarize the main findings. We also provide recommendations for data
230 scientists and statisticians, be it for industry, research or public organizations.

231 5.1 Limits

232 The obvious major drawback of the kNNxKDE is that we do not provide a clear way to optimize it.
233 We showed that our method performs best in terms of likelihood, but real-world datasets do not come
234 with a likelihood. Therefore, we are left with two options: either we use visual inspection and plots
235 to assess the data imputation quality, or we optimize τ to minimizing the RMSE (c.f. Appendix A).

236 Also, the kNNxKDE algorithm may not be suited for highly dimensional datasets. Not only can
237 it become computationally expensive, but its performances shall also worsen. Indeed, because of
238 the curse of dimensionality, initially close observations may end up far apart if similar features
239 are unobserved. This effect becomes even more problematic in high missing rates settings: as we
240 work with missing rate patterns, observations with few observed features will have a small number
241 of potential donors. This problem can be mitigated if the dataset has many observations. As a
242 consequence, calibrating the kNNxKDE algorithm in high dimensions is particularly challenging.
243 Pairplots may be used to visually assess the imputation quality, but become inconvenient in high-
244 dimension settings. Also, pairplots only display pairwise correlations and may overlook higher order
245 structures (c.f. Appendix B).

246 5.2 Strengths

247 If minimizing the imputation RMSE is an intuitive strategy for tabular data imputation, it cannot
248 capture the complexity of multimodal datasets. In practice, given an incomplete observation, if two
249 different imputations are consistent with the rest of the observed dataset, we have no objective way of
250 choosing one over the other. The kNNxKDE offers to not choose between these two options instead
251 of averaging over them both. It returns a imputation sample that provides more information than a
252 single point estimate.

253 Unlike the k NN-Imputer which impute column after column, the kNNxKDE works with successive
254 missing patterns. This allows to generate imputed samples which are consistent with the whole
255 dataset. Since all missing features are imputed at the same time, this strategy cannot return anomalous
256 imputed samples.

257 5.3 Conclusion

258 The main motivation of this work was to design an algorithm capable of imputing missing features
259 of a dataset with several modes. Multimodality makes imputation ambiguous, as clearly distinct
260 values may still be valid imputations. In this respect, we decide to use the likelihood as a metric
261 of imputation quality, instead of the standard RMSE between ground truth and imputed samples.
262 The kNNxKDE method does not aggregate estimations. Instead, it returns imputation samples all
263 consistent with the observed dataset. If needed, minimizing the imputation RMSE is possible by
264 averaging over the imputation samples, although we discourage from straightforwardly doing so as it
265 may lead to inconsistent imputed observations (c.f. Appendix A).

266 Ultimately, this work advocates for a qualitative approach of data imputation, rather than the current
267 quantitative one. We believe that missing data imputation should be done carefully and meaningfully,
268 as it influences subsequent data analysis. We provide the kNNxKDE algorithm, and we suggest trying
269 it for practical tabular data imputation in various domains.

270 References

271 Dimitris Bertsimas, Colin Pawlowski, and Ying Daisy Zhuo. From predictive methods to missing
272 data imputation: An optimization approach. *Journal of Machine Learning Research*, 18, 2018.
273 ISSN 15337928.

274 John K. Dixon. Pattern recognition with partly missing data. *IEEE Transactions on Systems, Man
275 and Cybernetics*, 9, 1979. ISSN 21682909. doi: 10.1109/TSMC.1979.4310090.

- 276 Nicholas Frosst, Nicolas Papernot, and Geoffrey Hinton. Analyzing and improving representations
277 with the soft nearest neighbor loss. In *ICML2019*, volume 2019-June, 2019.
- 278 Allison Marie Horst, Alison Presmanes Hill, and Kristen B Gorman. *palmerpenguins: Palmer*
279 *Archipelago (Antarctica) penguin data*, 2020. URL [https://allisonhorst.github.io/](https://allisonhorst.github.io/palmerpenguins/)
280 [palmerpenguins/](https://allisonhorst.github.io/palmerpenguins/). R package version 0.1.0.
- 281 Anil Jadhav, Dhanya Pramod, and Krishnan Ramanathan. Comparison of performance of data
282 imputation methods for numeric dataset. *Applied Artificial Intelligence*, 33, 2019. ISSN 10876545.
283 doi: 10.1080/08839514.2019.1637138.
- 284 Sebastian Jäger, Arndt Allhorn, and Felix Bießmann. A benchmark for data imputation methods.
285 *Frontiers in Big Data*, 4, 2021. ISSN 2624909X. doi: 10.3389/fdata.2021.693674.
- 286 Richard Leibrandt and Stephan Günnemann. Making kernel density estimation robust towards
287 missing values in highly incomplete multivariate data without imputation. In *SIAM2018*, 2018.
288 doi: 10.1137/1.9781611975321.84.
- 289 Roderick J.A. Little and Donald B. Rubin. *Statistical analysis with missing data*. Wiley, 2014. doi:
290 10.1002/9781119013563.
- 291 Emanuel Parzen. On estimation of a probability density function and mode. *The Annals of Mathe-*
292 *matical Statistics*, 33, 1962. ISSN 0003-4851. doi: 10.1214/aoms/1177704472.
- 293 Jason Poulos and Rafael Valle. Missing data imputation for supervised learning. *Applied Artificial*
294 *Intelligence*, 32, 2018. ISSN 10876545. doi: 10.1080/08839514.2018.1448143.
- 295 Murray Rosenblatt. Remarks on some nonparametric estimates of a density function. *The Annals of*
296 *Mathematical Statistics*, 27, 1956. ISSN 0003-4851. doi: 10.1214/aoms/1177728190.
- 297 Daniel J. Stekhoven and Peter Bühlmann. Missforest-non-parametric missing value imputation for
298 mixed-type data. *Bioinformatics*, 28, 2012. ISSN 13674803. doi: 10.1093/bioinformatics/btr597.
- 299 D. M. Titterton and G. M. Mill. Kernel-based density estimates from incomplete data. *Journal*
300 *of the Royal Statistical Society: Series B (Methodological)*, 45, 1983. doi: 10.1111/j.2517-6161.
301 1983.tb01249.x.
- 302 Olga Troyanskaya, Michael Cantor, Gavin Sherlock, Pat Brown, Trevor Hastie, Robert Tibshirani,
303 David Botstein, and Russ B. Altman. Missing value estimation methods for dna microarrays.
304 *Bioinformatics*, 17, 2001. ISSN 13674803. doi: 10.1093/bioinformatics/17.6.520.
- 305 Stef van Buuren and Karin Groothuis-Oudshoorn. mice: Multivariate imputation by chained equations
306 in r. *Journal of Statistical Software*, 45, 2011. ISSN 15487660. doi: 10.18637/jss.v045.i03.
- 307 Katarzyna Woznica and Przemyslaw Biecek. Does imputation matter? benchmark for predictive
308 models. *Artemiss2020, ICML workshop*, 2020. doi: 10.48550/arXiv.2007.02837.
- 309 Jinsung Yoon, James Jordon, and Mihaela Van Der Schaar. Gain: Missing data imputation using
310 generative adversarial nets. *35th International Conference on Machine Learning, ICML 2018*, 13:
311 9042–9051, 2018.

312 **A Real-world dataset: minimizing the RMSE with kNNxKDE**

313 For practical purposes, one may remain interested in minimizing the RMSE between the imputed
314 sample and the ground truth. This appendix shows how to use the kNNxKDE to obtain similar
315 RMSE performances as standard data imputation methods. The imputation samples returned by
316 the kNNxKDE allow for many ways of performing a single imputation. Rather than sampling the
317 conditional distributions only once for imputation – like we did in Section 4 – we can compute
318 appropriate statistics to estimate the missing values. Here, we use the mean for imputation.

319 The hyperparameter τ of the kNNxKDE is tuned to minimize the imputation NRMSE when using
320 the mean for the imputation. We use the Penguins dataset [Horst et al., 2020]: 342 penguins with 4
321 features (beak length, beak depth, flipper length and body mass) organized in 3 classes. This dataset

Table 4: Mean and standard deviation of the NRMSE on the Penguins dataset with all imputation methods. Optimal hyperparameters (shown below each method name) are obtained to minimize the NRMSE. k NNxKDE(m) stands for imputation performed with the mean of the returned samples from the k NNxKDE.

k NN-Imputer 40 neighbors	MissForest 30 trees	MICE x	GAIN 1200 iterations	k NNxKDE default	k NNxKDE(m) $\tau = 15$
0.136 ± 0.008	0.147 ± 0.012	0.154 ± 0.008	0.186 ± 0.026	0.219 ± 0.014	0.140 ± 0.012

322 is similar to the famous iris dataset. Results are reported in Table 4, where hyperparameters are
 323 optimized to minimize the NRMSE.

324 As we can see, averaging over the conditional distributions leads to similar performances as with the
 325 standard k NN-Imputer. The difference is that we now tune the continuous hyperparameter τ , which
 326 defines how loose the neighborhood of each observation is, rather than the number of neighbors k for
 327 the standard k NN-Imputer.

328 Note that, while the resulting imputation minimizes the RMSE, this may not preserve the structure of
 329 the original dataset any longer. If the original dataset is multimodal, the imputed dataset can present
 330 inconsistent observations.

331 B Synthetic data in 3d: visualizing higher-order correlations

332 We generate a dataset in 3-dimensions using spherical coordinates. Pairplots cannot help visualizing
 333 beyond pairwise correlations. But some structures may involve higher-order dependencies which
 334 traditional data imputation algorithms do not capture. For example, Figure 7 compares the imputation
 335 of the 3-d synthetic dataset with the k NN-Imputer and with the k NNxKDE. Table 5 presents the
 336 NRMSE and the log-likelihood for each method.

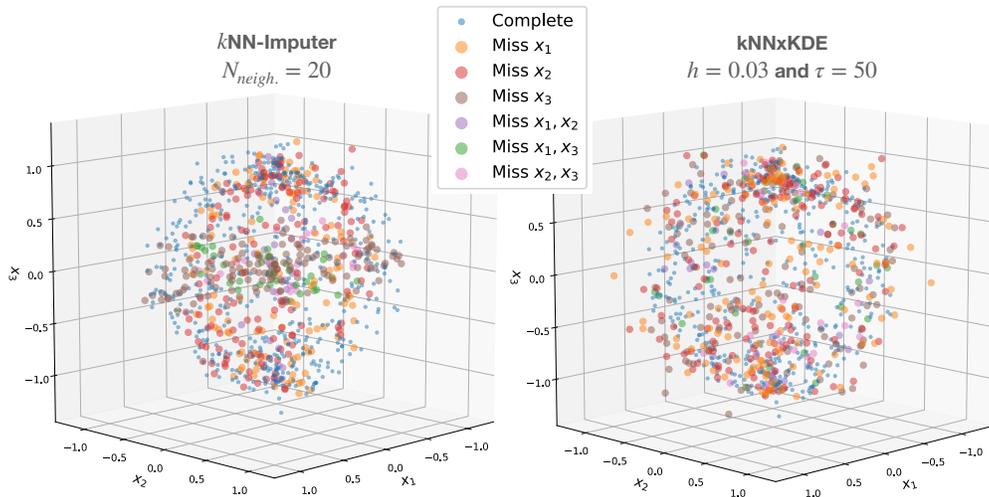


Figure 7: Visualization of the imputed 3-d spherical dataset (MCAR scenario with 20% missing rate): k NN-Imputer (left panel) and k NNxKDE (right panel). Points colors indicate imputed components. The k NN-Imputer creates artifacts (points inside the sphere) while the k NNxKDE preserve the original dataset structure.

337 Regarding the NRMSE, the k NNxKDE performs bad. But using the log-likelihood as benchmark, we
 338 see that the random sample generated by the k NNxKDE is much more probable under the generative
 339 model, i.e. the imputed sample is consistent with the original dataset. The scatter of the imputed
 340 observations (right panel of Figure 7) can be adjusted with τ and h .

341 Visual animations of the imputed samples with all five imputation methods are provided as supple-
 342 mentary materials, where we can notice the characteristics of each imputation method.

Table 5: Mean and standard deviation of the NRMSE on the Penguins dataset with all imputation methods. Optimal hyperparameters (shown below each method name) are obtained to minimize the NRMSE. kNNxKDE(m) stands for imputation performed with the mean of the returned samples from the kNNxKDE.

<i>(hyperparams)</i>	<i>k</i> NN-Imputer 20 neighbors	MissForest 15 trees	MICE x	GAIN 1200 iterations	kNNxKDE default
NRMSE	0.252	0.276	0.248	0.257	0.385
Log-Lik. (Ref=-2130)	-5683	-4023	-6309	-5793	-3008

343 Checklist

344 The checklist follows the references. Please read the checklist guidelines carefully for information on
 345 how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or
 346 **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing
 347 the appropriate section of your paper or providing a brief inline description. For example:

- 348 • Did you include the license to the code and datasets? **[Yes]** See Section xxx
- 349 • Did you include the license to the code and datasets? **[No]** The code and the data are
 350 proprietary.
- 351 • Did you include the license to the code and datasets? **[N/A]**

352 Please do not modify the questions and only use the provided macros for your answers. Note that the
 353 Checklist section does not count towards the page limit. In your paper, please delete this instructions
 354 block and only keep the Checklist section heading above along with the questions/answers below.

355 1. For all authors...

- 356 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 357 contributions and scope? **[Yes]** Emphasis on quality imputation and multimodal datasets
- 358 (b) Did you describe the limitations of your work? **[Yes]** See Section 5.1
- 359 (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
- 360 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 361 them? **[Yes]** To our knowledge, no potential negative or harmful societal impact. We
 362 have done our best for transparency and reproducibility

363 2. If you are including theoretical results...

- 364 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]**
- 365 (b) Did you include complete proofs of all theoretical results? **[N/A]**

366 3. If you ran experiments...

- 367 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 368 mental results (either in the supplemental material or as a URL)? **[Yes]** See supplement-
 369 ary materials
- 370 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 371 were chosen)? **[Yes]** Methodology and training procedures are extensively explained in
 372 Sections 2, 3 and 4
- 373 (c) Did you report error bars (e.g., with respect to the random seed after running exper-
 374 iments multiple times)? **[Yes]** Standard deviation are used for error bars (Section
 375 4). Seeds have been used in the code (supplementary materials) when needed for
 376 reproducibility
- 377 (d) Did you include the total amount of compute and the type of resources used (e.g.,
 378 type of GPUs, internal cluster, or cloud provider)? **[No]** We thought it was irrelevant,
 379 because rather fast with CPUs

380 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

- 381 (a) If your work uses existing assets, did you cite the creators? **[Yes]** We use one existing
 382 dataset, whose creators have been credited

- 383 (b) Did you mention the license of the assets? [Yes]
384 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
385 Code and synthetic data in Supplementary materials
386 (d) Did you discuss whether and how consent was obtained from people whose data you're
387 using/curating? [N/A]
388 (e) Did you discuss whether the data you are using/curating contains personally identifiable
389 information or offensive content? [N/A]
390 5. If you used crowdsourcing or conducted research with human subjects...
391 (a) Did you include the full text of instructions given to participants and screenshots, if
392 applicable? [N/A]
393 (b) Did you describe any potential participant risks, with links to Institutional Review
394 Board (IRB) approvals, if applicable? [N/A]
395 (c) Did you include the estimated hourly wage paid to participants and the total amount
396 spent on participant compensation? [N/A]