

Missing Data Imputation with Adversarially-trained Graph Convolutional Networks

Indro Spinelli^a, Simone Scardapane^{a,*}, Aurelio Uncini^a

^a*Department of Information Engineering, Electronics and Telecommunications (DIET), Sapienza University of Rome, Via Eudossiana 18, 00184 Rome, Italy*

Abstract

Missing data imputation (MDI) is the task of replacing missing values in a dataset with alternative, predicted ones. Because of the widespread presence of missing data, it is a fundamental problem in many scientific disciplines. Popular methods for MDI use global statistics computed from the entire dataset (e.g., the feature-wise medians), or build predictive models operating independently on every instance. In this paper we propose a more general framework for MDI, leveraging recent work in the field of graph neural networks (GNNs). We formulate the MDI task in terms of a graph denoising autoencoder, where each edge of the graph encodes the similarity between two patterns. A GNN encoder learns to build intermediate representations for each example by interleaving classical projection layers and locally combining information between neighbors, while another decoding GNN learns to reconstruct the full imputed dataset from this intermediate embedding. In order to speed-up training and improve the performance, we use a combination of multiple losses, including an adversarial loss implemented with the Wasserstein metric and a gradient penalty. We also explore a few extensions to the basic architecture involving the use of residual connections between layers, and of global statistics computed from the dataset to improve the accuracy. On a large experimental evaluation with varying levels of artificial noise, we show that our method is on par or better than several alternative imputation methods. On three datasets with pre-existing missing values, we show that our method is robust to the choice of a downstream classifier, obtaining similar or slightly higher results compared to other choices.

*Corresponding author. Phone: +39 06 44585495, Fax: +39 06 4873300.
Email address: `simone.scardapane@uniroma1.it` (Simone Scardapane)

Keywords: Imputation; Graph neural network; Graph data; Convolutional network

1. Introduction

While machine learning and deep learning have achieved tremendous results over the last years (Goodfellow et al., 2016), with new breakthroughs arising constantly (e.g., in drug discovery (Chen et al., 2018)), the vast majority of supervised learning methods still require datasets with complete information. At the same time, many real-world problems require dealing with incomplete data, such as in the biomedical or insurance sectors (Van Buuren, 2018). For this reason, flexible missing data imputation (MDI) methods are a fundamental component for widespread adoption of machine learning. An MDI algorithm takes a dataset with missing values in some of its input vectors, and replaces these values with some appropriately predicted ones in order to obtain a full dataset.¹ In particular, there is the need for powerful *multivariate* imputation methods able to work in a variety of data generation regimes (Yoon et al., 2018).

It has been recognized for a while that data imputation can be framed under a predictive framework, and classical machine learning methods (e.g., for regression and classification) might be adapted for this task (Bertsimas et al., 2017). However, some care must be taken when adapting them, for two fundamental reasons. Firstly, different inputs in general have different missing components, while most machine learning models assume full input vectors. Secondly, MDI might be a simple preprocessing step for downstream learning tasks, and in this case performance in terms of reconstruction might not be a perfect proxy for classification/regression accuracy later on.

In general, the resulting predictive approaches to MDI can be classified depending on whether they try to build a global model for data imputation, or whether they use similar data points to infer the missing components. Algorithms in the latter class include using simple statistics computed from the entire dataset (e.g., medians), or more advanced k-NN strategies (Lakshminarayan et al., 1996). In the former case, instead, we have simple linear

¹In the literature, this problem is also called missing value imputation (MVI) (Lin and Tsai, 2019). In this paper, we use the terms *data* and *value* interchangeably based on context.

models (Lakshminarayan et al., 1996), support vector machines (Wang et al., 2006) or, more recently, deep neural architectures (Yoon et al., 2018). These are surveyed more in-depth in Section 2.1.

We argue that a more powerful technique for MDI should exploit both ideas, i.e., use similar data points for each imputation *and* global models built from the overall dataset. In fact, recently a large class of neural network techniques have emerged that are able to model and exploit this kind of structured information (in the form of relationships between examples), by working in the domain of graphs (Bronstein et al., 2017; Battaglia et al., 2018). These models have been applied successfully to a wide range of problems, among which recommender systems (Ying et al., 2018), quantum chemistry (Gilmer et al., 2017), entity extraction from relational data (Schlichtkrull et al., 2018), semi-supervised learning (Kipf and Welling, 2017), and many others. However, to the best of our knowledge these techniques have never been applied to MDI. To overcome this, in this paper we define an architecture for MDI based on a specific class of graph neural networks, namely, graph convolutional networks (GCN), and empirically evaluate it on a large set of benchmark datasets (Kipf and Welling, 2017).

Contributions of the paper

Our generic framework for MDI is shown later on in Figure 1. We frame the overall problem in terms of a GCN autoencoder,² that learns to reconstruct the overall dataset conditioned on some artificial noise added during the training phase (similar to a classical denoising autoencoder (Vincent et al., 2008)). To build a graph of similarities between points we leverage prior literature on manifold regularization (Belkin et al., 2006), and we describe a simple technique that was found to work well in most situations.

After describing the basic architecture, we also detail three extensions to it that are able to improve either the accuracy or the speed of convergence:

- Firstly, we train the autoencoder with a mixture of standard loss functions and an adversarial loss, which was shown to provide significant improvements for denoising autoencoders in the non-graph case (Yoon et al., 2018).

²We use the term autoencoder to refer to any architecture that learns to map an input (or a corrupted version in the case of denoising autoencoders) to itself.

- Secondly, we motivate another extension with the inclusion of residual connections from the input to the output layer, similar to residual networks (He et al., 2016).
- Finally, we also describe how to include global information on the dataset (e.g., means and medians for all feature columns) using a generic context vector in input to the GCN layers.

We test our overall architecture on a large benchmark of datasets with varying levels of artificially-added noise and three real-world datasets with pre-existing missing values (two biomedical datasets and one time-series dataset). For the former, we show that our proposed GINN method is on par or outperforms several existing state-of-the-art approaches, especially when we consider high levels of injected artificial noise, e.g., up to 50% of missing values in the original dataset. For the latter, we show that our method is robust to the selection of a downstream classifier, with an accuracy comparable to any other combination of an imputation method and a classifier.

Organization of the paper

The rest of the paper is organized as follows. In Section 2 we describe the relation of this paper with state-of-the-art methods for MDI (Section 2.1) and graph neural networks (Section 2.2). The GCN, which is the building block of our method, is described in Section 3. Then, our graph imputation neural network (GINN) framework and all its extensions are described in Section 4. After a large experimental evaluation in Section 5, we provide some concluding remarks in Section 6.

2. Related work

2.1. Missing data imputation

Algorithms for MDI can be categorized depending on whether they perform univariate or multivariate imputation, and on whether they provide one or multiple imputations for each missing datum (Van Buuren, 2018). In addition, different algorithms can make different theoretical assumptions on whether the data is missing completely at random (MCAR) or not. In this paper we consider multivariate imputation, which is standard in the neural network’s literature. In the following we briefly review state-of-the-approaches on this topic, including several algorithms that we will compare to, and discuss their relation with our proposal.

A popular technique for MDI is multiple imputation using chained equations (MICE) (Azur et al., 2011; White et al., 2011; Van Buuren, 2018). MICE iteratively imputes each variable in the dataset by keeping the other variables fixed, repeating this for multiple cycles, each time drawing one or more observations from some predictive distribution on that variable. Although MICE has shown very good performance in some settings, especially in the bio-medical sector, the assumptions beyond MICE (especially the MCAR assumption) might result in biased predictions and subsequently lower accuracy (Azur et al., 2011).

In the machine learning community, it was recognized very soon that MDI can be framed as a predictive task, on which variants of standard supervised algorithms can be applied, including k-nearest neighbors (k-NN) (Acuna and Rodriguez, 2004), decision trees (Lakshminarayan et al., 1996), support vector techniques (Wang et al., 2006), and several others. However, these techniques have always achieved mixed performance in practice compared to simpler strategies such as mean imputation (Bertsimas et al., 2017). k-NN is limited in making weighted averages of similar feature vectors, while other algorithms are required to build a global model of the dataset to be used for imputation. In this paper we also frame the MDI problem in a predictive context, but our proposed model can leverage both global aspects of the dataset and local similarities between different points.

More recently, there has been a surge of interest in applying deep learning techniques to the problem of MDI. These include multiple imputation with deep denoising autoencoders (MIDA) (Gondara and Wang, 2018), combinations of deep networks with probabilistic mixture models (Śmieja et al., 2018), recurrent neural networks (Bengio and Gingras, 1996; Che et al., 2018), or generative models including generative adversarial networks (Yoon et al., 2018) and variational autoencoders (Nazabal et al., 2018). Generally speaking, these methods are better at capturing complex correlations in the data (and in the missing data process), thanks to their multiple layers of nonlinear computations, but they still require to build a global model from the dataset, while ignoring potentially important contributions from similar points. The method we propose can be seen as an extension both of the MIDA algorithm and of Yoon et al. (2018), but we focus on a more recent class of NNs, graph NNs, to capture local dependencies. We briefly survey the literature on this topic next.

2.2. Graph neural networks

Some of the earliest works on extending NNs to the domain of generic graphs were presented in [Gori et al. \(2005\)](#); [Scarselli et al. \(2009\)](#), and later reformulated in [Li et al. \(2015\)](#) in a more recent context. These works were mainly motivated by the analogies between unrolled recurrent neural networks and the diffusion of information across a graph.

Another line of work, upon which we build our proposal, considers instead the extension of convolutional neural networks to graph domains under the general term of geometric deep learning ([Defferrard et al., 2016](#); [Kipf and Welling, 2017](#); [Bronstein et al., 2017](#)) (and ([Micheli, 2009](#)) for earlier works on a similar context). This is done by exploiting recent ideas in the field of graph signal processing ([Sandryhaila and Moura, 2013](#); [Sardellitti et al., 2017](#)) to define a more general convolution operator able to work on irregular data structures. In particular, in this work we use the GCN of [Kipf and Welling \(2017\)](#), that for every layer includes a linear diffusion process across neighbors. Additional interesting lines of research in building GNNs that we briefly mention include earlier works on graph autoencoders ([Sperduti, 1994](#)), graph attention networks ([Veličković et al., 2018](#)), non-local NNs ([Wang et al., 2018](#)), graph embeddings ([Zhang et al., 2018](#)), and tree/graph echo state networks ([Gallicchio and Micheli, 2010, 2013](#)). An overview of many of these ideas is provided in [Battaglia et al. \(2018\)](#). We explore some of the ideas from [Battaglia et al. \(2018\)](#) in our framework by discussing how to include global information about the dataset in the reconstruction process in [Section 4.5](#).

Finally, our work is related to the field of manifold regularization ([Belkin et al., 2005, 2006](#)), a semi-supervised class of methods that exploits similarity information among patterns to enforce a regularization term on the optimization process. We build upon them for the construction of our similarity graph, a necessary step for exploiting the power of GNNs.

3. Graph convolutional networks

Because the GCN layer is a fundamental building block of our method, we briefly describe it here before moving on to the proposed framework for MDI. Consider a set of n vertices of a directed graph, whose connectivity is described by a (weighted) adjacency matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, where A_{ij} is different from 0 if and only if nodes i and j are connected. Each node i has an associated vector of features $\mathbf{x}_i \in \mathbb{R}^d$, that we collect row-wise in the matrix

$\mathbf{X} \in \mathbb{R}^{n \times d}$. We would like to have a generic neural network component able to process simultaneously the features at every node, but also take into consideration their relations, expressed through the adjacency matrix.

One way to extend the idea of convolutional networks to this domain is the so-called graph Fourier transform (Bruna et al., 2013; Sandryhaila and Moura, 2013). Define the Laplacian matrix of the graph as $\mathbf{L} = \mathbf{D} - \mathbf{A}$, where \mathbf{D} is the diagonal degree matrix with $D_{ii} = \sum_{j=1}^n A_{ij}$. We can perform the eigendecomposition of this matrix as $\mathbf{L} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$, where \mathbf{U} is a matrix collecting column-wise the eigenvectors of \mathbf{L} , and $\mathbf{\Lambda}$ is a diagonal matrix with the associated eigenvalues. The equivalent of a classical Fourier transform on a signal can be defined in the graph domain as (Sandryhaila and Moura, 2013):

$$\hat{\mathbf{X}} = \mathbf{U}^T \mathbf{X}, \quad (1)$$

and the inverse transform as $\mathbf{X} = \mathbf{U}\hat{\mathbf{X}}$. Using this, a straightforward way to define a convolutional layer on graphs (Bruna et al., 2013) is to first apply the graph Fourier transform, apply a trainable transformation on the frequency components (associated to the eigenvalues of the Laplacian), and then back-transform using the inverse Fourier transform. While viable, this approach is however costly and impractical in most cases.

Later authors (Defferrard et al., 2016; Kipf and Welling, 2017) have noted that by applying a restricted class of filters to the frequency components (polynomials), it is possible to work directly in the graph domain using polynomials of the Laplacian itself. In particular, Kipf and Welling (2017) proposed the GCN with the use of linear filters, resulting in the following canonical layer:

$$\mathbf{H}_1 = g(\mathbf{L}\mathbf{X}\mathbf{\Theta}_1), \quad (2)$$

where $\mathbf{\Theta}_1$ is a matrix of adaptable coefficients, and $g(\cdot)$ a generic element-wise activation function, such as the ReLU $g(s) = \max(0, s)$.³ Note that the right-multiplication by $\mathbf{\Theta}_1$ is akin to a classical feedforward layer, while the left-multiplication by \mathbf{L} allows to propagate the information across the immediate neighbors of each node. Multiple layers of this form can then be

³To avoid some numerical instabilities, it is possible to renormalize the Laplacian to properly bound its eigenvalues, as done in Kipf and Welling (2017). More in general, one can substitute the Laplacian with any valid graph shift operator (Gama et al., 2019).

stacked to obtain a complete graph NN. Importantly, for a generic network with L layers of the form (2), the output of node i will depend on the outputs of its neighbors up to degree L .

4. Proposed framework for missing data imputation

In MDI, we are also given a data matrix \mathbf{X} , which has the same size and semantic as in the previous section, but in general no graph information associated with it. Some of the values of \mathbf{X} , denoted by a binary mask $\mathbf{M} \in \{0, 1\}^{n \times d}$, are missing and need to be imputed for downstream processing or classification/regression. We assume to have either numerical features, which are properly normalized, or categorical features that are represented with one-hot encoding. For other types of features, e.g., text, a previous embedding step is needed (Pennington et al., 2014).⁴

Predictive models for MDI described previously in Section 2.1 build a function $f(\mathbf{x}_i)$ for imputing missing values of a single example \mathbf{x}_i , but in general, do not exploit directly the potentially important information contained in points that might be similar to it. Here, we propose to model this constraint explicitly by building f using GCN blocks, as shown schematically in Figure 1. To do this, we first need to build a graph of inter-patterns similarities from \mathbf{X} , as described in the next section.

4.1. Construction of the similarity graph

The first fundamental step of our method is the discovery of the graph structure underneath the tabular representation of the data. In the resulting graph, each feature vector in the dataset is encoded as a node of the graph, while the adjacency matrix \mathbf{A} is derived from a similarity matrix \mathbf{S} of the features vectors. As we stated before, constructing a similarity graph from the data is a known problem in the literature, and here we leverage some work from the field of manifold regularization (Belkin et al., 2005, 2006), adapting it for handling the presence of missing data. A small overview on possible alternatives is provided at the end of this section.

⁴When facing a supervised learning problem, for which there is an additional label (e.g., class) associated to each input \mathbf{x}_i , we can easily include the training labels in the imputation process by concatenating them to the input vector.

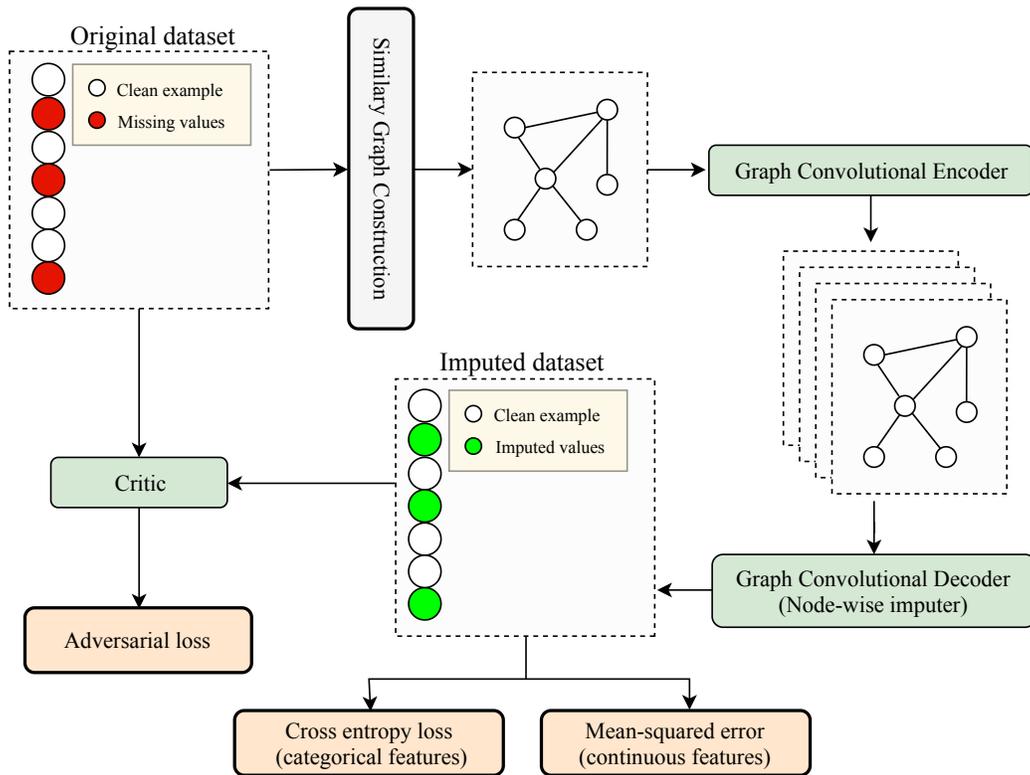


Figure 1: Schematics of the proposed framework for missing data imputation. In green we show the components that are trained end-to-end on the imputation task. In orange we show the different losses. Details on all the steps are provided in the text.

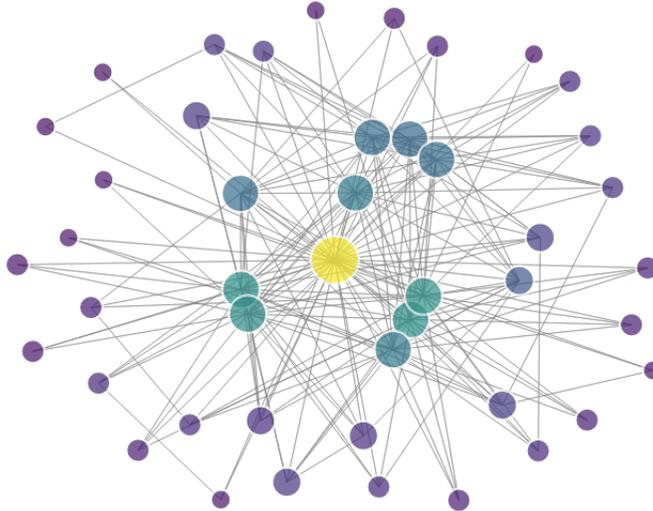


Figure 2: Subset of the graph reconstructed from the Iris dataset using our pipeline. For each node the color intensity represents the number of connections and the size the number of missing features. A clear correlation can be seen between the two.

The similarity matrix is computed pairwise for all features vectors using the Euclidean distance, but each time only the non-missing elements of both vectors are used for the computation (Eirola et al., 2013):

$$S_{ij} = d(\mathbf{x}_i \odot (\mathbf{M}_i \odot \mathbf{M}_j), \mathbf{x}_j \odot (\mathbf{M}_i \odot \mathbf{M}_j)) \quad (3)$$

where \odot stands for the Hadamard product between vectors, \mathbf{M}_i is the i th column of the binary matrix \mathbf{M} , and d is the Euclidean distance.

In practice, it is common to have sparse similarity matrices (Belkin et al., 2006), most notably for efficiency and computational reasons (see in particular the discussion on computational cost later on), allowing most operations to scale at-most linearly in the number of neighbors. In order to have a sparse graph, with meaningful connections between similar nodes, we apply a pruning step over the similarity matrix \mathbf{S} inspired from the large-scale manifold learning algorithm in Talwalkar et al. (2013). A threshold is applied independently over every row of \mathbf{S} , by computing a percentile for each row that will act as a threshold, such that only the connections above this threshold inside every row are kept. To make this process more robust, we iterate it twice. The result is used as adjacency matrix for the GCN in the next section.

We found that the 97.72nd percentile provides good results over all the datasets used in the benchmark of Section 5, allowing us to discard at each step around 95% of all the possible connections (in accordance with Talwalkar et al. (2013)). Figure 2 shows a subset of the graph produced with this procedure starting from the Iris dataset; here we display the relationship between the number of missing features in a node and its degree. It can be seen from the figure how nodes with very few non-zero elements are correlated with a higher number of connections.⁵

On the construction of the similarity graph: the method described in this section, which is the one we follow in our experiments and in our open-source implementation, was found to provide good empirical performance. Nonetheless, we underline that in the proposed GINN framework, the similarity graph can be built according to any guideline or method available in the literature. For example, classical alternative choices in the semi-supervised literature include binary weights on the edges, heat kernel similarity (Belkin et al., 2006), or selecting a fixed number of neighbors instead of a percentile (Geng et al., 2012). If the inputs contain text, images, or similar data, cosine similarity on custom pre-trained embeddings are also a popular choice (Bui et al., 2018). We leave an analysis of these different alternatives to future work.

4.2. Autoencoder architecture

Autoencoders are composed by an encoder which maps the inputs to an intermediate representation in a different dimensional space $\mathbf{h} = \text{encode}(\mathbf{x})$, and a decoder that maps $\mathbf{h} \in \mathbb{R}^m$ to the original dimensional space $\hat{\mathbf{x}} = \text{decode}(\mathbf{h})$. We use $m > d$ for an overcomplete representation, thus mapping the input into a higher dimensional space with the aim of helping data recovery. Our graph imputer neural network (GINN) will thus be defined as follows:

$$\begin{aligned} \mathbf{H} &= \text{ReLU}(\mathbf{LX}\Theta_1) \\ \hat{\mathbf{X}} &= \text{Sigmoid}(\mathbf{LH}\Theta_2) \end{aligned} \tag{4}$$

⁵This is due to the absence of a re-normalization step in the computation of the similarity matrix (Eirola et al., 2013). In practice, we have found this setup to work better than renormalizing all distance measures, possibly because of the increased degree of elements with multiple missing values.

where \mathbf{L} has been defined in Section 3 (the extension to networks with multiple hidden layers being straightforward).

Note that we cannot trivially train the autoencoder on the missing values, because they are not known in the training stage. To solve this, we adopt a denoising version of the autoencoder (Vincent et al., 2008), in which for every optimization step we add additional masking noise on the input, by the means of an inverted dropout layer applied directly on the input of the network. In particular, at each optimization step, we randomly remove 50% of the remaining inputs.⁶ In this way, the autoencoder learns to reconstruct any part of the input matrix, similarly to (Gondara and Wang, 2018).

We train the whole model end-to-end minimizing the reconstruction error over the non-missing elements of the dataset. The loss function is thus defined as the combination of a mean squared error (MSE) for the numerical variables and the cross-entropy (CE) for the categorical variables:

$$L_A = \alpha \text{MSE}(\mathbf{X}, \hat{\mathbf{X}}) + (1 - \alpha) \text{CE}(\mathbf{X}, \hat{\mathbf{X}}) \quad (5)$$

where MSE always returns 0 for categorical values and *vice versa* for CE. α is an additional hyper-parameter that we initialize as the ratio between the number of numerical columns of the dataset and the total number of columns. Alternatively, it can be tuned like the other hyper-parameters of the network, although we have not found any definite improvement in doing so.

Computational cost of the model

The computational cost of our approach is related mostly to (a) the one-time cost of constructing the similarity graph, and (b) the use of a GCN layer instead of a standard feedforward layer as in alternative autoencoder architectures (Gondara and Wang, 2018). The cost of point (a) is well-studied in the literature on large-scale similarity search, e.g., (Dong et al., 2011), and many techniques and implementations exist for making it efficient. The cost of the GCN layer is discussed in Kipf and Welling (2017, Section 3.2). In particular, using a sparse representation for the adjacency matrix reduces the memory requirement to $\mathcal{O}(|E|)$, where E is the number of edges in the graph, and the cost of computing (2) to $\mathcal{O}(|E|CF)$, where C and F

⁶The only exception: whenever training labels are used as inputs for the imputation process, we do not apply dropout on them.

are the number of input and output features respectively. In practice, when using an early-stopping strategy for training, we have found the training time for our GINN algorithm (including point (a)) to be significantly faster than alternative neural approaches and on-par with non-neural competitors such as MICE, e.g., see Fig. 4 in the additional materials for a comparison.

4.3. Adversarial training of the autoencoder

In order to speed up training, we use an additional adversarial training strategy where a critic, a feedforward network in our case, learns to distinguish between imputed and real data. This is inspired by generative adversarial networks (Goodfellow et al., 2014) and was found to have significant effects in several reconstructions tasks, particularly in the medical domain (Ker et al., 2018). In particular, having an adversarial loss during reconstruction forces the reconstructed vector to lie close to the natural distribution of the original patterns (Shen et al., 2019).

To train jointly autoencoder and critic and have a stable training we used the Wasserstein distance introduced in Arjovsky et al. (2017), which is informally defined as the minimum cost of transporting mass in order to transform a distribution q into a distribution p . Using the Kantorovich-Rubinstein duality (Villani, 2008) the objective function is obtained as follows:

$$\min_A \max_{C \in \mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{real}} [C(\mathbf{x})] - \mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{imp}} [C(\hat{\mathbf{x}})] \quad (6)$$

where \mathcal{D} is the set of 1-Lipschitz functions, \mathbb{P}_{imp} is the model distribution implicitly defined by our GCN autoencoder $\hat{\mathbf{x}} = A(\mathbf{x})$, and \mathbb{P}_{real} is the unknown data distribution. Practically, Eq. (6) can be computed by drawing random mini-batches of data and approximating expectations with averages.

The original loss in Arjovsky et al. (2017) used weight clipping to force the Lipschitz property. A further step towards training stability, introduced in Gulrajani et al. (2017), is to use a gradient penalty instead of the weight clipping, obtaining the final loss:

$$L_C = \mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{imp}} [C(\hat{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{real}} [C(\mathbf{x})] + \lambda \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_{\tilde{\mathbf{x}}}} [(\|\nabla_{\tilde{\mathbf{x}}} C(\tilde{\mathbf{x}})\|_2 - 1)^2] \quad (7)$$

where λ is an additional hyper-parameter. We define $\mathbb{P}_{\tilde{\mathbf{x}}}$ as sampling uniformly from the combination of the real distribution \mathbb{P}_{imp} and from the distribution resulting from the imputation \mathbb{P}_{imp} . This means that the feature

vector $\tilde{\mathbf{x}}$ will be composed by both real and imputed elements in almost equal quantity. This distance is continuous and differentiable thus, the more we train the critic, the more reliable the gradients are. Following standard practice in the GAN literature, in our implementation the GCN autoencoder is trained once for every five optimization steps of the critic, and the total loss becomes:

$$L_D = L_A - \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_{imp}} [C(\hat{\mathbf{x}})] \quad (8)$$

since it must fool the critic and minimize the reconstruction error at the same time.

4.4. Including skip connections in the model

The autoencoder itself generates an approximate reconstruction of the dataset, while the critic loss guides the autoencoder in this process. However, our main scope is the imputation of values not present in the data. For this task we want a greater contribution coming from the most similar nodes. For this reason, we introduce an additional skip layer which consists always in a graph convolution operation but propagating the information across the immediate neighbors of each node without the node itself. This prevents the autoencoder from learning the identity function.

The decoding layer becomes:

$$\hat{\mathbf{X}} = \text{Sigmoid} \left(\mathbf{LH}\Theta_2 + \tilde{\mathbf{L}}\mathbf{X}\Theta_3 \right) \quad (9)$$

where $\tilde{\mathbf{L}}$ is computed similarly to \mathbf{L} (as described in Section 3), but starting from an adjacency matrix without self loops. This is shown schematically in Figure 3.

4.5. Including global statistics from the dataset

Another extension we explore is the possibility of including global information on the dataset during the computation of the autoencoder. The inclusion of a global set of attributes, in the context of graph neural networks, was described in-depth by Battaglia et al. (2018).

As a proof of concept, in our case we set as global attribute vector \mathbf{g} for the graph some statistical information of the dataset, including mean or mode of every attribute. The global component can be taken into account in the last layer, weighting their contribution for the update of each node:

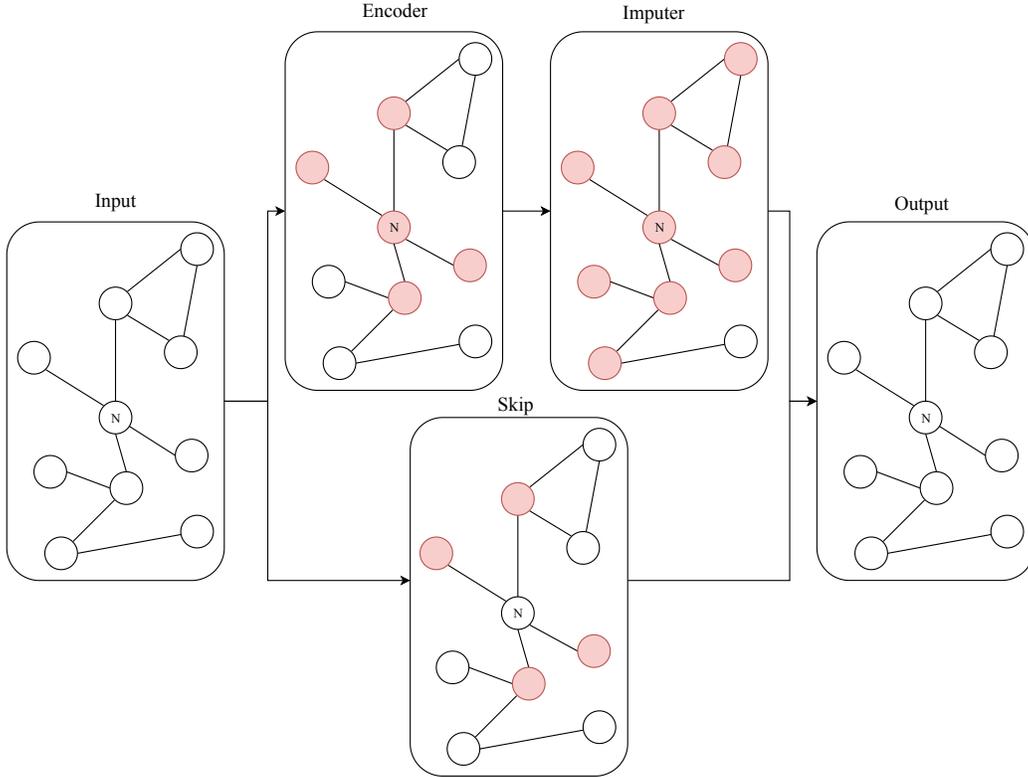


Figure 3: Schematics of the graph autoencoder with the skip connection. We highlight the nodes involved, directly or indirectly by the convolution for the update of the node N . The encoding phase involves 1-hop neighbors, decoding/imputation phase instead involves nodes up to 2-hops. The skip connections involves 1-hop neighbors without considering node N .

$$\hat{\mathbf{X}} = \text{Sigmoid} \left(\mathbf{LH}\Theta_2 + \tilde{\mathbf{L}}\mathbf{X}\Theta_3 + \Theta_4\mathbf{g} \right) \quad (10)$$

In addition, if the computation of the global information is differentiable, we can compute a loss term with respect to the global attributes of the original dataset:

$$L = L_d + \gamma \text{MSE}(\text{global}(\hat{\mathbf{X}}), \text{global}(\mathbf{X})) \quad (11)$$

where γ is some additional weighting term.

5. Experimental evaluation

We divide our experimental evaluation in five subsections. Firstly, following common literature, we evaluate the proposed GINN framework on 20 real-world datasets from the UCI Machine Learning Repository (Dua and Graff, 2017) to which we artificially add some desired level of missing values, to evaluate the imputation performances. The characteristics of these 20 datasets are summarized in Table 1. This selection contains categorical, numerical, and mixed datasets, ranging from 150 observation to 30000 and from only 4 attributes to almost 40. Every dataset is divided into training 70% and test 30% sets. Missingness is introduced completely at random on the training set with 4 different levels of noise: 10%, 20%, 30%, and 50%. Our evaluation will focus first on imputation performance as described in Section 5.1, then on the accuracy of post-imputation prediction in Section 5.2. We then perform a comprehensive ablation study of the architecture in Section 5.3, and an evaluation of the performance of the algorithm on new data in Section 5.4.

Secondly, in Section 5.5 we evaluate the performance of the algorithm on three real-world datasets with pre-existing (i.e., non artificially induced) missing values. In this part we also evaluate the computational cost of the method when compared to other state-of-the-art approaches. In one case, missing values are also present in the training labels, in which case we provide an application of our method to a semi-supervised scenario (as described later on).

For all the benchmarks, we use an embedding dimension of the hidden layer of 128, sufficient for an overcomplete representation for all the datasets involved, and we train the model for a maximum of 10000 iterations with an early stopping strategy for the reconstruction loss over the known elements. The critic used is a simple 3-layer feed-forward network trained 5 times for each optimization step of the autoencoder. We used the Adam optimizer (Kingma and Ba, 2014) for both networks with a learning rate of 1×10^{-3} and 1×10^{-5} respectively for autoencoder and critic. When label information is available for the datasets, we consider the training labels as an additional feature of each input vector, but we remove this information when processing new (validation or test) data. All experiments are repeated five times and we collect average performance and standard deviation.

All the code for replicating our experiments and using the GINN algo-

rithm is released as an open-source library on the web.⁷

Name	observations	numerical attr.	categorical attr.
abalone	4177	8	0
anuran-calls	7195	22	3
balance-scale	625	4	0
breast-cancer-diagnostic	569	30	0
car-evaluation	1728	0	6
default-credit-card	30000	13	10
electrical-grid-stability	10000	14	0
heart	303	8	5
ionosphere	351	34	0
iris	150	5	0
page-blocks	5473	10	0
phishing	1353	0	9
satellite	6435	36	0
tic-tac-toe	958	0	9
turkiye-student-evaluation	5820	0	32
wine-quality-red	1599	11	0
wine-quality-white	4898	11	0
wine	178	13	0
wireless-localization	2000	7	0
yeast	1484	8	0

Table 1: Datasets used for the benchmark. All of them were downloaded from the UCI repository.

5.1. Imputation Performance

This evaluation focuses on the comparison of MAE and RMSE between GINN and 6 other state-of-the-art imputation algorithms: MICE ([van Buuren and Groothuis-Oudshoorn, 2011](#)), MIDA ([Gondara and Wang, 2018](#)), MissForest ([Stekhoven and Buehlmann, 2012](#)), mean ([Little and Rubin, 1986](#)), matrix factorization and k-NN imputation ([Botstein et al., 2001](#)). Concerning MICE and MissForest we used the default parameters discussed in the corresponding papers. For the other methods, we fine-tuned the hyperparameters according to the corresponding literature to provide a fair com-

⁷<https://github.com/spindro/GINN>

parison. In particular, we used 1×10^{-3} and 1×10^{-4} as learning rate for matrix factorization and MIDA with the latter being a 2-layer with 128 units per layer network like our embedding dimension. Finally, we let $k = 5$ for the k-NN. The imputation accuracy for each dataset is presented in Table 2 for the scenario in which 30% of the entries are missing, while the results for all other levels of missingness (both in terms of RMSE and MAE) are presented in the supplementary material. We can see from Table 2 that the proposed GINN method obtains the best imputation performance in almost half of the datasets, being the second-best in almost all the remaining ones.

To provide a more schematic comparison, in Figure 4 we show the summary of those results for every level of missing data. In these histograms, we provide the number of times that each method achieves the best imputation (on average), i.e. the lowest RMSE in Figure 4(a) and MAE in Figure 4(b). Concerning the lower percentage of missing elements (10%, 20%) our method is almost on par with the best among the algorithms tested, i.e., MissForest. When those percentages increase our method brings a huge improvement over the state-of-the-art with the highest difference at 50% where our method significantly outperforms all other techniques. Aggregating the results obtained at 30% and 50% percentage of missing features, we have that our method is the best in 50% of the cases against the 27.5% of its best competitor MissForest, when looking at the MAE, and 47.5% against 20% for the RMSE. We defer a statistical analysis of these results to the next subsection, where we analyze also the results for a downstream predictive task.

5.2. Predictive Performance

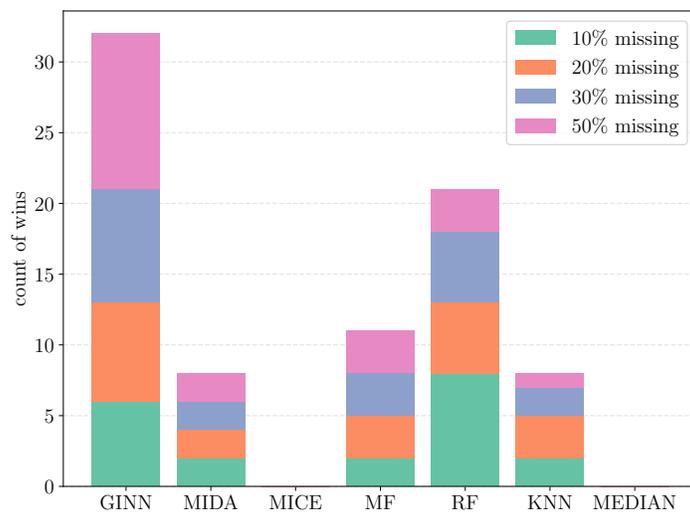
Now we evaluate the performance of standard machine learning algorithms for classification, both binary and multi-class, trained on the various imputations analyzed previously. We consider 4 different classifiers: a k-NN classifier with $k = 5$, regularized logistic regression, C-Support Vector Classification with an RBF kernel and a random forest classifier with 10 estimators and a maximum depth of 5. All hyper-parameters are initialized with their default values in the scikit-learn implementation.⁸

The classification accuracy is presented in Table 3 for the scenario in which 30% of the entries in the data matrix are missing, assuming MCAR,

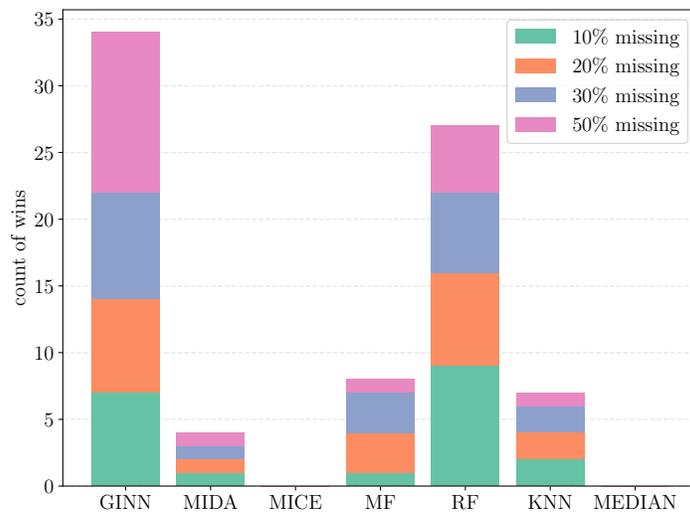
⁸<https://scikit-learn.org/stable/modules/classes.html>

RMSE	GINN	MIDA	MICE	MF	RF	k-NN	MEDIAN
abalone	0.904 ± 0.048	1.190 ± 0.002	1.069 ± 0.053	3.321 ± 2.220	0.811 ± 0.008	1.072 ± 0.035	1.112 ± 0.001
anuran-calls	0.070 ± 0.009	0.187 ± 0.001	0.076 ± 0.000	0.091 ± 0.000	0.154 ± 0.003	0.058 ± 0.000	0.226 ± 0.002
balance-scale	0.559 ± 0.010	0.436 ± 0.001	0.579 ± 0.005	0.521 ± 0.021	0.516 ± 0.000	0.580 ± 0.002	0.577 ± 0.001
breast-cancer-diagnostic	54.633 ± 4.519	114.390 ± 5.612	38.122 ± 3.568	23.559 ± 2.964	21.155 ± 2.401	39.653 ± 6.124	127.645 ± 6.907
car-evaluation	0.623 ± 0.003	0.632 ± 0.003	0.636 ± 0.003	0.845 ± 0.000	0.636 ± 0.002	0.636 ± 0.006	0.649 ± 0.002
default-credit-card	17479.409 ± 280.991	19664.115 ± 197.618	15687.056 ± 262.499	15212.927 ± 388.170	12897.406 ± 279.296	17294.412 ± 56.008	23708.027 ± 100.798
electrical-grid-stability	1.534 ± 0.005	1.857 ± 0.003	1.637 ± 0.005	1.474 ± 0.023	1.536 ± 0.008	1.763 ± 0.028	1.555 ± 0.001
heart	10.883 ± 2.003	25.212 ± 1.926	12.326 ± 1.506	11.367 ± 1.745	11.368 ± 1.814	14.846 ± 1.360	11.708 ± 1.726
ionosphere	0.425 ± 0.001	1.118 ± 0.005	0.425 ± 0.007	340.257 ± 20.027	0.407 ± 0.006	0.380 ± 0.001	0.551 ± 0.004
iris	0.349 ± 0.047	1.792 ± 0.106	0.379 ± 0.003	0.106 ± 0.269	0.333 ± 0.008	0.446 ± 0.062	1.159 ± 0.053
page-blocks	536.645 ± 44.594	1035.750 ± 80.323	1889.068 ± 6.830	795.441 ± 85.233	579.227 ± 147.154	604.171 ± 33.633	869.476 ± 49.517
plishing	0.493 ± 0.006	0.529 ± 0.006	0.625 ± 0.000	0.594 ± 0.000	0.606 ± 0.020	0.548 ± 0.001	0.581 ± 0.002
satellite	6.256 ± 0.200	11.528 ± 0.112	4.481 ± 0.013	4.562 ± 0.024	3.634 ± 0.005	4.523 ± 0.040	18.335 ± 0.060
tic-tac-toe	0.469 ± 0.010	0.576 ± 0.000	0.657 ± 0.001	0.816 ± 0.005	0.671 ± 0.005	0.625 ± 0.001	0.622 ± 0.001
turkiye-student-evaluation	0.295 ± 0.002	0.285 ± 0.000	0.537 ± 0.001	0.287 ± 0.000	0.292 ± 0.001	0.303 ± 0.000	0.515 ± 0.001
wine	47.901 ± 7.298	154.344 ± 22.959	52.358 ± 14.401	54.764 ± 6.209	48.672 ± 8.529	54.636 ± 8.874	99.651 ± 10.909
wine-quality-red	8.099 ± 0.097	10.313 ± 0.091	9.092 ± 0.079	8.883 ± 0.100	8.414 ± 0.071	9.752 ± 0.079	10.265 ± 0.031
wine-quality-white	11.188 ± 1.133	15.383 ± 0.351	11.174 ± 0.183	41.324 ± 0.196	10.350 ± 0.153	12.454 ± 0.391	13.432 ± 0.223
wireless-localization	4.851 ± 0.227	8.470 ± 0.131	4.728 ± 0.047	1.863 ± 0.113	4.214 ± 0.105	4.995 ± 0.010	8.588 ± 0.129
yeast	0.085 ± 0.000	0.162 ± 0.002	0.090 ± 0.000	0.093 ± 0.000	0.086 ± 0.000	0.091 ± 0.001	0.102 ± 0.002

Table 2: Average Root Mean Squared Error with 30% of missing elements. Best results for each dataset is highlighted with a bold font. One standard deviation is also provided.



(a) Comparison over MAE



(b) Comparison over RMSE

Figure 4: Number of datasets in which each MDI method achieves (a) lowest average MAE or (b) lowest average RMSE from the true values. The different colors of the bar stand for different percentages of missing elements: from the bottom 10% at the lowest to the top 50% at the highest.

with a random forest classifier. In Figure 5 we show the summary of the results for each noise level and for each classifier. In these histograms we

analyze the number of times each imputation technique allows the classifier to achieve the best average accuracy. In this comparison, we consider also the draws.

Our method outperforms competitors with every classification algorithm tested, and it has the highest number of wins, winning in 85.62% of the cases. Our method worked very well when paired with SVC and Random forest, improving the accuracy for every percentage of missing elements. With the logistic regression and k-NN classifiers, at the lowest missing percentage (10%), our method is slightly below the state-of-the-art. As missing percentages increase, we have a huge improvement over the other competitors. This reflects in the findings of the previous Section and confirms the ability of our method of being very successful in the case of moderate to severely damaged datasets.

We corroborate the results of this and the previous section by performing a statistical analysis of the algorithms according to the guidelines in [Demšar \(2006\)](#). A Friedman rank test confirms that there are statistical significant differences both with respect to the RMSE of Figure 4 (p-value of $1.11e^{-11}$), and with respect to the accuracy of the classifiers in Figure 5 (e.g., p-value for random forest is $1.01e^{-10}$). A successive set of Nemenyi post-hoc tests further confirms statistical significant differences between GINN and all other methods for the random forest classification in Figure 5, and between GINN and all other methods except RF for the imputation results in Figure 4. The full set of rankings and of p-values for these tests can be found in the additional material for this paper.

5.3. Ablation study

To investigate how much each step described in Section 4 improves our method, we supervised the quality of imputation and convergence by starting from a basic autoencoder. As the starting point we used a 2-layer denoising autoencoder (DAE), obtained by setting the adjacency matrix to be the identity, obtaining a method similar to ([Gondara and Wang, 2018](#)). Then we introduced the graph and the graph convolutional layer in our imputer (GINN) followed by the addition of the critic and the adversarial training (A-GINN), the skip connection (A-GINN skip) and finally the global attributes (A-GINN skip global).

Regarding imputation performances, Table 4 shows how the introduction of the graph and the graph-convolution operation over three randomly selected datasets makes a huge difference against a standard autoencoder.

Random forest	baseline	GINN	MIDA	MICE	MF	RF	k-NN	MEDIAN
abalone	53.269 ± 0.804	55.821 ± 1.569	53.269 ± 1.169	54.944 ± 1.037	51.754 ± 1.1675	52.472 ± 1.2036	53.030 ± 0.040	51.674 ± 1.555
anuran-calls	92.728 ± 0.463	94.441 ± 0.787	89.810 ± 1.760	91.848 ± 0.648	93.191 ± 0.556	92.635 ± 0.394	91.060 ± 0.162	90.412 ± 1.135
balance-scale	81.328 ± 0.0687	76.595 ± 0.452	72.340 ± 0.982	76.595 ± 1.053	77.659 ± 1.301	76.063 ± 0.834	78.723 ± 0.893	69.148 ± 1.140
breast-cancer-diagnostic	97.076 ± 0.877	96.491 ± 0.585	95.906 ± 2.636	94.736 ± 0.877	94.152 ± 1.170	97.076 ± 0.877	96.491 ± 1.170	97.660 ± 1.170
car-evaluation	70.520 ± 0.674	72.832 ± 0.586	71.483 ± 0.096	69.942 ± 0.096	69.942 ± 0.000	70.712 ± 0.393	70.520 ± 0.482	69.942 ± 0.000
default-credit-card	77.870 ± 0.100	77.980 ± 0.001	77.880 ± 0.000	77.880 ± 0.000	77.886 ± 0.042	77.893 ± 0.017	77.9 ± 0.025	77.9 ± 0.083
electrical-grid-stability	98.533 ± 0.033	91.100 ± 1.633	86.634 ± 0.867	97.556 ± 1.650	94.990 ± 3.967	97.867 ± 0.000	98.700 ± 0.333	99.360 ± 1.233
heart	82.417 ± 2.747	83.516 ± 0.549	73.626 ± 2.198	83.516 ± 1.648	78.021 ± 2.198	82.417 ± 1.648	82.417 ± 0.549	76.923 ± 3.297
ionosphere	90.566 ± 0.001	95.283 ± 2.358	84.905 ± 1.415	92.452 ± 0.000	92.452 ± 1.415	92.452 ± 0.943	95.283 ± 1.415	93.396 ± 1.415
iris	91.111 ± 0.000	88.889 ± 1.667	80.000 ± 1.667	93.334 ± 0.000	91.112 ± 3.333	91.112 ± 3.333	93.334 ± 1.667	86.667 ± 0.000
page-blocks	94.640 ± 0.594	95.066 ± 0.137	94.823 ± 0.000	94.762 ± 0.137	94.640 ± 0.046	94.640 ± 0.046	94.884 ± 0.137	94.701 ± 0.000
phishing	84.729 ± 1.107	84.482 ± 0.738	83.004 ± 0.738	83.251 ± 0.369	83.990 ± 0.923	83.004 ± 0.738	84.236 ± 0.554	82.019 ± 0.185
satellite	83.990 ± 0.207	83.850 ± 1.157	81.699 ± 0.155	82.550 ± 0.181	83.650 ± 0.958	83.300 ± 0.233	83.500 ± 0.363	79.950 ± 0.984
tic-tac-toe	69.443 ± 0.868	72.569 ± 0.174	68.056 ± 0.694	67.708 ± 0.868	67.361 ± 1.042	65.972 ± 4.688	71.180 ± 1.910	72.916 ± 0.174
turkiye-student-evaluation	84.020 ± 0.200	84.650 ± 0.859	84.593 ± 0.143	81.214 ± 0.601	84.650 ± 0.830	83.276 ± 0.086	83.505 ± 1.088	83.218 ± 0.286
wine	96.296 ± 0.000	98.148 ± 0.000	92.592 ± 1.852	98.148 ± 0.926	94.444 ± 0.926	98.148 ± 0.000	94.444 ± 0.926	94.444 ± 3.704
wine-quality-red	63.334 ± 1.655	62.708 ± 2.105	62.083 ± 3.635	61.667 ± 2.105	63.334 ± 1.766	62.291 ± 1.042	61.458 ± 1.766	58.125 ± 1.655
wine-quality-white	52.108 ± 0.544	53.673 ± 0.952	51.360 ± 0.578	52.517 ± 0.374	52.312 ± 1.224	53.265 ± 0.204	51.020 ± 0.136	51.836 ± 0.204
wireless-localization	97.166 ± 0.167	96.334 ± 1.083	93.167 ± 1.417	97.334 ± 0.833	96.667 ± 0.333	97.000 ± 0.167	95.834 ± 0.083	97.167 ± 0.083
yeast	55.156 ± 0.561	56.502 ± 1.457	46.188 ± 3.027	58.071 ± 0.785	43.721 ± 0.897	54.484 ± 0.112	56.502 ± 2.018	53.363 ± 2.018

Table 3: Classification accuracy on each dataset using a Random forest classifier. The model was trained over the imputed data by algorithms for 30% MCAR. In the first column (baseline) we have the accuracy obtained utilizing the undamaged dataset.

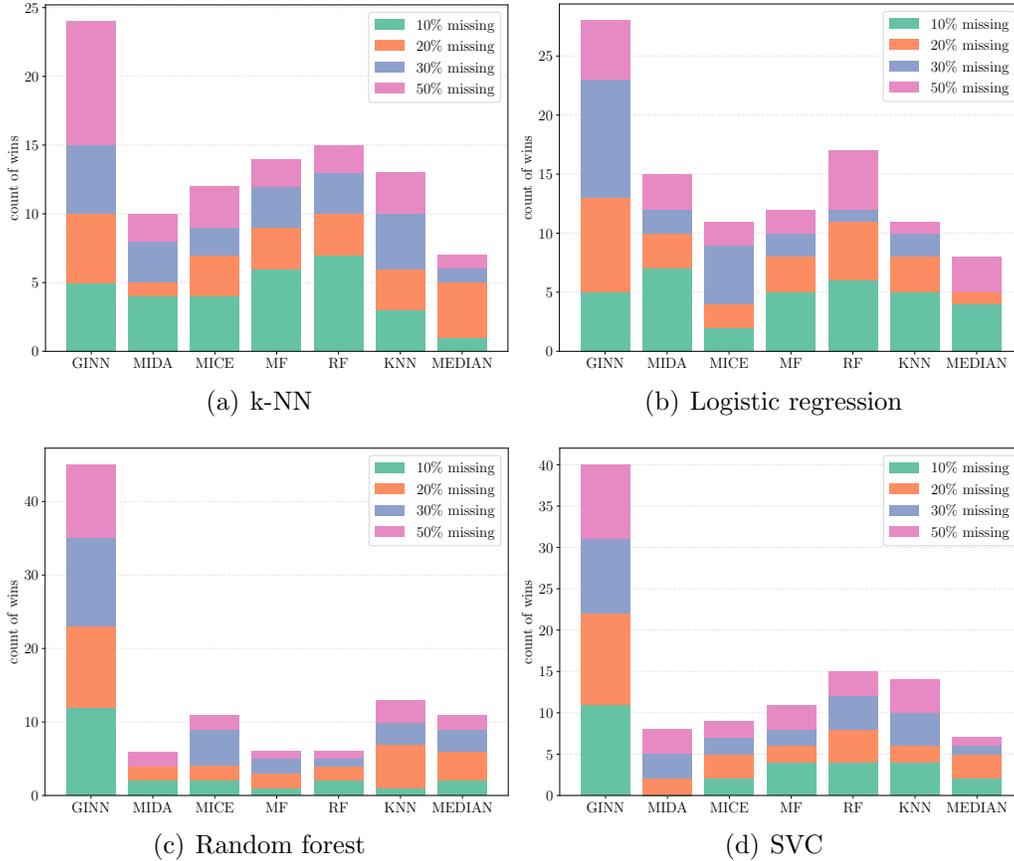
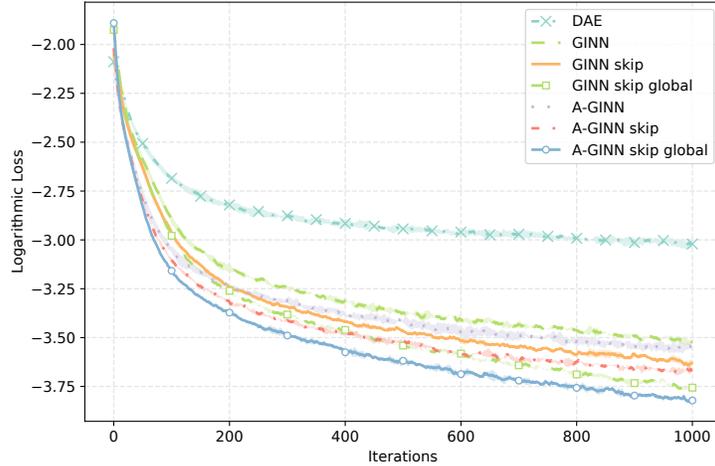


Figure 5: Number of datasets in which each MDI method achieves the highest classification accuracy for each classifier used. The different colors of the bar stand for different percentages of missing elements: from the bottom 10% the lowest to the top 50% the highest.

After that, each following step helps to refine even further the imputation accuracy. The reconstruction loss in Eq. (5), more precisely its logarithm, shows a similar behaviour, shown in Figure 6. After each step we have a better convergence. Similar results are obtained for all the other datasets.

The results in Figure 6 are shown with respect to the number of iterations. In the supplementary material, we also provide a similar analysis with respect to a fixed computational budget, while an analysis of the overall computational time when compared to the other algorithms is provided later on in Section 5.5.



(a)

Figure 6: Convergence of the logarithmic loss function defined in Equation 5, on the Ionosphere dataset for each improvement described in this section and with imputation results shown in Table 4.

5.4. Imputation over unseen data

In this section we test the ability of the model of imputing a new damaged portion of the dataset that was not available at training time. In order to impute these new values, we have first to inject the new data in the graph, adding nodes and edges. We compute a new similarity matrix for the new features (not considering the labels) and also the similarity between these new features and the older ones. Then we add the new nodes to the graph and the edges resulting from the double threshold procedure described in Section 4.1.

To evaluate this, we introduced missing values and evaluated the imputation performances with and without fine-tuning the model on a second randomly kept portion of the datasets. In Table 5 we show the MAE of the imputation over this new portion of dataset and compare against the other MDI algorithms. We used the same dataset and settings of the ablation study. The fine-tuned version consists of an additional 500 epochs of training over the new graph. It can be seen how our method is able to perform a state-of-the-art imputation on new unseen data without performing additional training and how it improves in case of a small number of additional optimization steps.

	Ionosphere	Tic-Tac-Toe	Phishing
DAE	0.309 ± 0.011	0.323 ± 0.004	0.260 ± 0.005
GINN	0.263 ± 0.013	0.317 ± 0.003	0.247 ± 0.002
GINN skip	0.258 ± 0.013	0.314 ± 0.002	0.246 ± 0.007
GINN skip global	0.256 ± 0.013	0.313 ± 0.000	0.245 ± 0.007
A-GINN	0.257 ± 0.011	0.316 ± 0.001	0.243 ± 0.006
A-GINN skip	0.256 ± 0.013	0.305 ± 0.007	0.243 ± 0.006
A-GINN skip global	0.255 ± 0.016	0.303 ± 0.003	0.241 ± 0.004

Table 4: Mean Absolute Error, \pm the standard deviation, of the imputation over 3 trials with 20% of missing elements. Each dataset has the corresponding convergence shown in Figure 6.

MAE	Ionosphere	Tic-Tac-Toe	Phishing
GINN	<u>0.252 ± 0.008</u>	<u>0.320 ± 0.012</u>	<u>0.242 ± 0.002</u>
GINN (FT)	0.235 ± 0.004	0.299 ± 0.005	0.237 ± 0.002
MIDA	0.782 ± 0.000	0.390 ± 0.004	0.422 ± 0.006
MICE	0.335 ± 0.008	0.414 ± 0.004	0.644 ± 0.003
MF	0.312 ± 0.009	0.667 ± 0.003	0.462 ± 0.029
RF	0.254 ± 0.024	0.456 ± 0.001	0.344 ± 0.000
k-NN	0.235 ± 0.002	0.400 ± 0.006	0.282 ± 0.005
MEDIAN	0.362 ± 0.004	0.388 ± 0.004	0.343 ± 0.008

Table 5: Mean Absolute Error of the imputation with 20% of missing elements.

5.5. Evaluation on datasets with pre-existing missing values

To evaluate GINN in a real-world scenario, we compared its performance against the other state-of-the-art algorithms over three datasets with *pre-existing* missing values. When the data is not missing completely at random, the problem gets more complex, because there may exist relationships between the probability of a variable to be missing and other observed data. We consider two biomedical datasets: mammographic mass introduced in [Elter et al. \(2007\)](#) and cervical cancer by [Fernandes et al. \(2017\)](#) respectively with 4% and 13% of missing elements. We performed the imputation without considering the additional information of the label. Then we solved the binary classification task as done in Section 5.2. The third dataset is a time-series of air quality measurements ([De Vito et al., 2008](#)) with 13% of missing elements. Differently from the other two, missing values can also be found

Name	observations	numerical attr.	categorical attr.	missingness
c. cancer	860	9	25	13%
m. mass	960	0	5	4%
air quality	3313	9	0	13%

Table 6: Datasets with pre-existing missingness used for the benchmark in Section 5.5. All of them were downloaded from the UCI repository.

in the labels, making this a semi-supervised task. We select the three initial most damaged months as training data and perform imputation on all data, including the target variable. Then, we performed the downstream task of classifying, in the two successive months, the target variable discretized in three bins. A summary of the three datasets is provided in Table 6.

In Table 7 we show the average accuracy (computed over 10 trials) for every combination of imputation method and downstream classifier. For clarity, we highlight in bold the best result, and we underline the second-best one. While our proposed algorithm does not necessarily achieve the best accuracy overall, it can be seen from Table 7 that it is, in average, the most resilient to the choice of an external classifier. Overall, when combined with a logistic regression we obtain the best accuracy for the cervical cancer dataset (on par with several other algorithms), while we obtain the second-best result when combined with a random forest or an SVC in the other two cases. In order to highlight the resilience of the algorithm, in the supplementary material we provide a statistical analysis when the results from the different classifiers are aggregated.

Concerning computational performance, we provide execution times for all the algorithms (for simplicity, in the random forest case) in the supplementary material. Briefly, our algorithm is faster than alternative neural approaches (as already described in Section 5.3), but slower than non-neural approaches such as k-NN.

6. Conclusions and future work

In this paper we introduced a novel technique for missing data imputation, where we used a novel graph convolutional autoencoder to reconstruct the full dataset. We also describe several improvement to our technique, including the use of an adversarial loss, and the inclusion of global information from the dataset in the reconstruction phase. We show through an extensive numerical simulation that our method has good imputation performance,

	Classifier	M. Mass	C. Cancer	A. Quality
GINN	k-NN	82.45±0.8	98.35±0.6	90.09±0.6
	LR	82.81±1.6	99.44±0.4	90.72±0.1
	RF	83.04±1.2	98.82±0.5	<u>90.86±0.5</u>
	SVC	<u>83.98±1.3</u>	97.71±0.1	90.08±0.0
MIDA	k-NN	81.18±0.7	98.00±0.6	88.29±0.6
	LR	81.70±1.5	<u>98.97±0.6</u>	89.53±0.1
	RF	78.50±2.8	97.77±0.2	88.99±0.9
	SVC	83.67±1.4	97.71±0.1	88.49±0.0
MICE	k-NN	80.45±1.2	97.83±0.4	80.99±0.8
	LR	82.41±1.2	99.44±0.4	81.58±0.0
	RF	80.19±1.9	97.83±0.4	78.14±0.9
	SVC	83.75±1.1	97.71±0.1	78.56±0.1
MF	k-NN	80.57±1.2	98.45±0.7	84.41±0.5
	LR	82.27±1.1	99.44±0.4	84.72±1.2
	RF	77.66±3.6	97.78±0.2	83.78±1.4
	SVC	84.03±1.2	97.71±0.1	83.29±0.0
RF	k-NN	81.02±1.4	97.78±0.2	88.41±0.4
	LR	82.67±1.1	99.44±0.4	88.30±1.7
	RF	76.66±6.8	97.89±0.5	88.65±1.1
	SVC	83.71±1.0	97.71±0.1	88.36±0.0
KNN	k-NN	81.61±1.2	98.10±0.6	89.38±0.5
	LR	82.37±1.8	99.44±0.4	91.18±0.0
	RF	80.11±2.0	97.81±0.3	89.24±0.5
	SVC	84.03±1.2	97.71±0.1	89.79±0.0
MEDIAN	k-NN	80.37±1.0	98.06±0.7	88.85±0.4
	LR	82.17±1.1	99.44±0.4	88.46±0.0
	RF	76.74±5.8	97.98±0.3	89.83±0.5
	SVC	84.03±1.2	97.71±0.1	88.36±0.0

Table 7: Classification accuracy and standard deviation over 10 trials obtained by using the different imputation algorithms on datasets having pre-existing missing values. For each dataset, we highlight in bold the best result, and we underline the second-best one.

and the results are robust to the selection of an additional classifier later on. In experiments with a large level of artificial noise, our method is also shown to significantly outperform competitors.

Future work can consider the adoption of different graph neural architectures for the autoencoding process (such as those mentioned in Section 2), or the extension to other types of noisy data beyond vector-valued data and different types of similarity measures. In addition, in order to further improve accuracy and training time, we can think of training our imputation module together with a classification step in an end-to-end fashion.

Currently, the major drawbacks of our method are the need for computing the similarity matrix of the data, and the difficulty of performing mini-batching in the presence of graph-based data. Both problems are well-known in the corresponding literature, and in future work we plan on investigating techniques for speeding up similarity search and mini-batching on the graph to improve the computational complexity of the method.

References

References

- Acuna, E., Rodriguez, C., 2004. The treatment of missing values and its effect on classifier accuracy. In: Classification, clustering, and data mining applications. Springer, pp. 639–647.
- Arjovsky, M., Chintala, S., Bottou, L., 2017. Wasserstein generative adversarial networks. In: Proc. 34th International Conference on Machine Learning (ICML). Vol. 70. pp. 214–223.
- Azur, M. J., Stuart, E. A., Frangakis, C., Leaf, P. J., 2011. Multiple imputation by chained equations: what is it and how does it work? *International Journal of Methods in Psychiatric Research* 20 (1), 40–49.
- Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., et al., 2018. Relational inductive biases, deep learning, and graph networks. arXiv preprint arXiv:1806.01261.
- Belkin, M., Niyogi, P., Sindhvani, V., 2005. On manifold regularization. In: AISTATS.

- Belkin, M., Niyogi, P., Sindhwani, V., 2006. Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *Journal of Machine Learning Research* 7 (Nov), 2399–2434.
- Bengio, Y., Gingras, F., 1996. Recurrent neural networks for missing or asynchronous data. In: *Advances in Neural Information Processing Systems*. pp. 395–401.
- Bertsimas, D., Pawlowski, C., Zhuo, Y. D., 2017. From predictive methods to missing data imputation: An optimization approach. *Journal of Machine Learning Research* 18, 196–1.
- Botstein, D., Sherlock, G., Cantor, M., Troyanskaya, O., Brown, P., Tibshirani, R., Altman, R. B., Hastie, T., 06 2001. Missing value estimation methods for DNA microarrays . *Bioinformatics* 17 (6), 520–525.
- Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., Vandergheynst, P., 2017. Geometric deep learning: going beyond euclidean data. *IEEE Signal Processing Magazine* 34 (4), 18–42.
- Bruna, J., Zaremba, W., Szlam, A., LeCun, Y., 2013. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*.
- Bui, T. D., Ravi, S., Ramavajjala, V., 2018. Neural graph learning: Training neural networks using graphs. In: *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*. pp. 64–71.
- Che, Z., Purushotham, S., Cho, K., Sontag, D., Liu, Y., 2018. Recurrent neural networks for multivariate time series with missing values. *Scientific Reports* 8 (1), 6085.
- Chen, H., Engkvist, O., Wang, Y., Olivecrona, M., Blaschke, T., 2018. The rise of deep learning in drug discovery. *Drug Discovery Today* 23 (6), 1241–1250.
- De Vito, S., Massera, E., Piga, M., Martinotto, L., 02 2008. On field calibration of an electronic nose for benzene estimation in an urban pollution monitoring scenario. *Sensors and Actuators B Chemical* 129, 750–757.
- Defferrard, M., Bresson, X., Vandergheynst, P., 2016. Convolutional neural networks on graphs with fast localized spectral filtering. In: *Advances in Neural Information Processing Systems*. pp. 3844–3852.

- Demšar, J., 2006. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research* 7 (Jan), 1–30.
- Dong, W., Moses, C., Li, K., 2011. Efficient k-nearest neighbor graph construction for generic similarity measures. In: *Proceedings of the 20th international conference on World wide web*. ACM, pp. 577–586.
- Dua, D., Graff, C., 2017. UCI machine learning repository.
URL <http://archive.ics.uci.edu/ml>
- Eirola, E., Doquire, G., Verleysen, M., Lendasse, A., 2013. Distance estimation in numerical data sets with missing values. *Information Sciences* 240, 115–128.
- Elter, M., Schulz-Wendtland, R., Wittenberg, T., 11 2007. The prediction of breast cancer biopsy outcomes using two cad approaches that both emphasize an intelligible decision process. *Medical physics* 34, 4164–72.
- Fernandes, K., Cardoso, J., Fernandez, J., 05 2017. Transfer learning with partial observability applied to cervical cancer screening. pp. 243–250.
- Gallicchio, C., Micheli, A., 2010. Graph echo state networks. In: *Proc. 2010 IEEE International Joint Conference on Neural Networks (IJCNN)*. IEEE, pp. 1–8.
- Gallicchio, C., Micheli, A., 2013. Tree echo state networks. *Neurocomputing* 101, 319–337.
- Gama, F., Marques, A. G., Leus, G., Ribeiro, A., 2019. Convolutional neural network architectures for signals supported on graphs. *IEEE Transactions on Signal Processing* 67 (4), 1034–1049.
- Geng, B., Tao, D., Xu, C., Yang, L., Hua, X.-S., 2012. Ensemble manifold regularization. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 34 (6), 1227–1233.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., Dahl, G. E., 2017. Neural message passing for quantum chemistry. In: *Proc. 34th International Conference on Machine Learning (ICML)*. JMLR. org, pp. 1263–1272.

- Gondara, L., Wang, K., 2018. Multiple imputation using deep denoising autoencoders. In: Proc. 22nd Pacific-Asia Conference on Knowledge Discovery and Data Mining (PAKDD). pp. 1–12.
- Goodfellow, I., Bengio, Y., Courville, A., 2016. Deep learning. MIT Press.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets. In: Advances in Neural Information Processing Systems. pp. 2672–2680.
- Gori, M., Monfardini, G., Scarselli, F., 2005. A new model for learning in graph domains. In: Proc. 2005 IEEE International Joint Conference on Neural Networks (IJCNN). Vol. 2. IEEE, pp. 729–734.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., Courville, A. C., 2017. Improved training of wasserstein gans. In: Advances in Neural Information Processing Systems. pp. 5767–5777.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. In: Proceedings 2016 IEEE Conference on Computer Vision and Pattern Recognition (ICCVPR). pp. 770–778.
- Ker, J., Wang, L., Rao, J., Lim, T., 2018. Deep learning applications in medical image analysis. IEEE Access 6, 9375–9389.
- Kingma, D. P., Ba, J., 2014. Adam: A method for stochastic optimization. In: Proc. 3rd International Conference for Learning Representations (ICLR).
- Kipf, T. N., Welling, M., 2017. Semi-supervised classification with graph convolutional networks. In: Proc. 2017 International Conference on Learning Representations (ICLR).
- Lakshminarayan, K., Harp, S. A., Goldman, R. P., Samad, T., et al., 1996. Imputation of missing data using machine learning techniques. In: Proc. KDD-96. pp. 140–145.
- Li, Y., Tarlow, D., Brockschmidt, M., Zemel, R., 2015. Gated graph sequence neural networks. In: Proc. 2016 International Conference on Learning Representations (ICLR). pp. 1–20.
- Lin, W.-C., Tsai, C.-F., 2019. Missing value imputation: a review and analysis of the literature (2006–2017). Artificial Intelligence Review, 1–23.

- Little, R. J. A., Rubin, D. B., 1986. *Statistical Analysis with Missing Data*. John Wiley & Sons, Inc., New York, NY, USA.
- Micheli, A., 2009. Neural network for graphs: A contextual constructive approach. *IEEE Transactions on Neural Networks* 20 (3), 498–511.
- Nazabal, A., Olmos, P. M., Ghahramani, Z., Valera, I., 2018. Handling incomplete heterogeneous data using vaes. arXiv preprint arXiv:1807.03653.
- Pennington, J., Socher, R., Manning, C., 2014. Glove: Global vectors for word representation. In: *Proc. 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. pp. 1532–1543.
- Sandryhaila, A., Moura, J. M., 2013. Discrete signal processing on graphs. *IEEE Transactions on Signal Processing* 61 (7), 1644–1656.
- Sardellitti, S., Barbarossa, S., Di Lorenzo, P., 2017. On the graph fourier transform for directed graphs. *IEEE Journal of Selected Topics in Signal Processing* 11 (6), 796–811.
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., Monfardini, G., 2009. The graph neural network model. *IEEE Transactions on Neural Networks* 20 (1), 61–80.
- Schlichtkrull, M., Kipf, T. N., Bloem, P., Van Den Berg, R., Titov, I., Welling, M., 2018. Modeling relational data with graph convolutional networks. In: *European Semantic Web Conference*. Springer, pp. 593–607.
- Shen, G., Dwivedi, K., Majima, K., Horikawa, T., Kamitani, Y., 2019. End-to-end deep image reconstruction from human brain activity. *Frontiers in Computational Neuroscience* 13.
- Śmieja, M., Struski, Ł., Tabor, J., Zieliński, B., Spurek, P., 2018. Processing of missing data by neural networks. In: *Advances in Neural Information Processing Systems*. pp. 2724–2734.
- Sperduti, A., 1994. Encoding labeled graphs by labeling raam. In: *Advances in Neural Information Processing Systems*. pp. 1125–1132.
- Stekhoven, D. J., Buehlmann, P., 2012. Missforest - non-parametric missing value imputation for mixed-type data. *Bioinformatics* 28 (1), 112–118.

- Talwalkar, A., Kumar, S., Mohri, M., Rowley, H., 2013. Large-scale svd and manifold learning. *The Journal of Machine Learning Research* 14 (1), 3129–3152.
- Van Buuren, S., 2018. Flexible imputation of missing data. Chapman and Hall/CRC.
- van Buuren, S., Groothuis-Oudshoorn, K., 2011. mice: Multivariate imputation by chained equations in r. *Journal of Statistical Software* 45 (3), 1–67.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y., 2018. Graph attention networks. In: *Proc. 2018 International Conference on Learning Representations (ICLR)*.
- Villani, C., 2008. Optimal transport – Old and new. Springer Science & Business Media.
- Vincent, P., Larochelle, H., Bengio, Y., Manzagol, P.-A., 2008. Extracting and composing robust features with denoising autoencoders. In: *Proc. 25th International Conference on Machine Learning (ICML)*. ICML '08. ACM, New York, NY, USA, pp. 1096–1103.
- Wang, X., Girshick, R., Gupta, A., He, K., 2018. Non-local neural networks. In: *Proc. IEEE Conference on Computer Vision and Pattern Recognition (ICCVPR)*. pp. 7794–7803.
- Wang, X., Li, A., Jiang, Z., Feng, H., 2006. Missing value estimation for dna microarray gene expression data by support vector regression imputation and orthogonal coding scheme. *BMC Bioinformatics* 7 (1), 32.
- White, I. R., Royston, P., Wood, A. M., 2011. Multiple imputation using chained equations: issues and guidance for practice. *Statistics in Medicine* 30 (4), 377–399.
- Ying, R., He, R., Chen, K., Eksombatchai, P., Hamilton, W. L., Leskovec, J., 2018. Graph convolutional neural networks for web-scale recommender systems. In: *Proc. 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. ACM, pp. 974–983.

- Yoon, J., Jordon, J., van der Schaar, M., 2018. Gain: Missing data imputation using generative adversarial nets. In: Proc. 35th International Conference of Machine Learning (ICML). pp. 1–10.
- Zhang, H., Wang, S., Xu, X., Chow, T. W., Wu, Q. J., 2018. Tree2vector: learning a vectorial representation for tree-structured data. *IEEE Transactions on Neural Networks and Learning Systems* 28 (11), 5304–5318.

Additional material

Detailed RMSE results (Section 5.1)

In Table 8 we provide the individual RMSE imputation values, where each row is one of the 20 datasets and each column is an imputation method (abbreviations are explained in the main text). We separate the results with respect to the level of artificial corruption of the original dataset: the suffix *_xx* means *xx%* of missing values that have been artificially added. These results are aggregated and commented in Fig. 4 and Tab. 3 of the main text. Results for MAE are similar and can be found on our online repository.

Detailed accuracy for regression/classification (Section 5.2)

In Table 9 we report the accuracy of the downstream classification/regression task for each choice of classification/regression technique when using a random forest technique (results are similar for other methods, and for brevity we provide them on our online repository). Like before, we separate the results with respect to the level of artificial corruption of the original dataset: the suffix *_xx* means *xx%* of missing values that have been artificially added. These results are aggregated and commented in Fig. 5(c) and Tab. 4 of the main text.

Detailed results for the statistical tests

In Fig. 7 we provide the average rankings and p-values for the statistical test performed on the results of Section 5.1 of the paper (Nemenyi post-hoc tests on all pairs of algorithms), corresponding to Tab. 2. In Fig. 8 we instead provide average rankings and corresponding p-values for the statistical tests performed on the random forest classifier of Section 5.2, corresponding to Tab. 9. The results are discussed more in-depth in the main paper in Section 5.2.

Ablation study with a computational budget (Section 5.3)

In this section we replicate the ablation study of Section 5.3, but we evaluate the convergence of each variant with respect to a fixed computational budget, as shown in Fig. 9. As can be seen, our baseline method with GCN but no adversarial loss can converge to a significantly better result than a standard DAE in a fraction of the time. Including the adversarial loss slightly improves the final result, requiring however 3 – 4 times the budget of the baseline method.

Detailed results for the evaluation on real-world datasets

In Figure 10 we report the scores obtained with a random forest classifier and the time in seconds needed for the imputation, including for GINN, the time needed for the similarity-graph construction. As can be seen, GINN’s imputation allows to achieve the best accuracy and consistency across all three datasets, even in the case where the imputed training labels are used to train the classification model. Regarding execution times, GINN is considerably faster than alternative neural approaches and missForest, but it is slower than the approaches that do not perform a training phase.

We corroborate the results by performing a statistical analysis. To focus on the impact of the imputation techniques, we compute the relative ranking for each combination of imputation method and classification algorithm, and we then aggregate the results with respect to the latter. A Friedman rank test confirms that there are statistical significant differences with respect to the accuracy of the classifiers with p-values of $1.6e^{-7}$, $1.3e^{-7}$, $1.2e^{-9}$ respectively for the Cervical cancer, Mammographic mass and Air quality datasets. A successive set of Nemenyi post-hoc tests, between all pairs of algorithms, further confirms statistical significant differences between GINN and all other methods as reported in Figure 11.

RMSE	GINN	MIDA	MICE	MF	RF	KNN	MEDIAN
abalone_10	0.901	0.987	1.002	3.110	0.693	0.741	1.031
abalone_20	0.919	1.087	1.031	3.335	0.783	1.031	1.164
abalone_30	0.904	1.190	1.069	3.321	0.811	1.072	1.112
abalone_50	1.038	1.218	1.052	2.679	0.948	1.121	1.155
anuran-calls_10	0.066	0.156	0.069	0.745	0.159	0.049	0.225
anuran-calls_20	0.067	0.157	0.073	0.750	0.158	0.052	0.221
anuran-calls_30	0.070	0.187	0.076	0.091	0.154	0.058	0.226
anuran-calls_50	0.081	0.202	0.086	0.097	0.132	0.089	0.227
balance-scale_10	0.565	0.424	0.579	0.545	0.552	0.614	0.591
balance-scale_20	0.575	0.426	0.559	0.540	0.558	0.604	0.586
balance-scale_30	0.559	0.436	0.579	0.521	0.516	0.580	0.577
balance-scale_50	0.548	0.441	0.566	0.495	0.482	0.581	0.577
breast-cancer_10	45.901	65.698	35.837	25.006	6.021	28.299	105.975
breast-cancer_20	39.258	84.307	35.248	23.413	19.332	27.415	122.510
breast-cancer_30	54.633	114.390	38.122	23.559	21.155	39.653	127.645
breast-cancer_50	57.873	136.427	37.325	26.060	28.434	48.219	121.975
car-evaluation_10	0.616	0.621	0.645	0.846	0.631	0.660	0.653
car-evaluation_20	0.613	0.619	0.628	0.845	0.637	0.647	0.646
car-evaluation_30	0.623	0.632	0.636	0.845	0.636	0.636	0.649
car-evaluation_50	0.615	0.629	0.641	0.846	0.636	0.648	0.640
credit-card_10	17585.4	18323.0	16034.1	15863.5	11762.4	15939.1	24680.3
credit-card_20	17792.4	18462.5	15848.6	14864.0	12055.4	15907.4	23884.7
credit-card_30	17479.4	19664.1	15687.0	15212.9	12897.4	17294.4	23708.0
credit-card_50	19437.2	20063.7	16394.4	15471.8	14139.7	20974.4	23652.4
electrical-grid_10	1.493	1.871	1.625	1.375	1.398	1.586	1.574
electrical-grid_20	1.494	1.855	1.638	1.318	1.438	1.624	1.538
electrical-grid_30	1.534	1.857	1.637	1.474	1.536	1.763	1.555
electrical-grid_50	1.570	1.878	1.656	1.686	1.724	1.666	1.580
heart_10	9.125	17.100	9.696	9.901	9.869	10.477	9.553
heart_20	8.068	18.842	11.220	11.363	10.046	12.816	10.900
heart_30	10.883	25.212	12.326	11.367	11.368	14.846	11.708
heart_50	10.301	20.874	11.671	10.743	10.910	11.651	11.004
ionosphere_10	0.364	0.772	0.416	0.572	0.349	0.348	0.527
ionosphere_20	0.385	0.819	0.458	351.366	0.386	0.379	0.542
ionosphere_30	0.425	1.118	0.462	340.257	0.407	0.380	0.551
ionosphere_50	0.442	1.170	0.475	185.546	0.431	0.415	0.544
iris_10	0.400	1.988	0.386	0.116	0.237	0.287	0.895
iris_20	0.384	1.562	0.411	0.133	0.310	0.352	1.261
iris_30	0.349	1.792	0.379	0.106	0.333	0.446	1.159
iris_50	0.385	1.992	0.453	0.109	0.534	0.414	1.135
page-blocks_10	407.7	814.3	1511.9	684.5	206.7	241.6	1019.8
page-blocks_20	780.0	2305.8	2006.4	1349.1	934.1	723.1	2512.4
page-blocks_30	536.6	1035.7	1889.0	795.4	579.2	604.1	869.4
page-blocks_50	1438.4	1873.5	1624.9	1044.6	1291.9	1858.2	1801.9

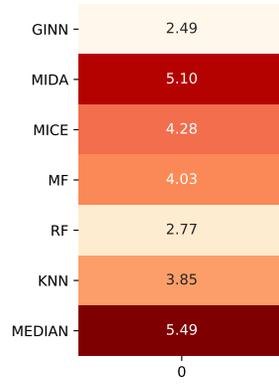
RMSE	GINN	MIDA	MICE	MF	RF	KNN	MEDIAN
phishing_10	0.500	0.506	0.618	0.520	0.607	0.519	0.576
phishing_20	0.513	0.514	0.623	0.531	0.601	0.524	0.579
phishing_30	0.493	0.529	0.625	0.594	0.606	0.548	0.581
phishing_50	0.506	0.541	0.623	0.581	0.596	0.538	0.576
satellite_10	5.978	9.251	4.162	4.846	3.207	4.077	18.293
satellite_20	6.066	9.802	4.312	4.158	3.404	4.257	18.423
satellite_30	6.256	11.528	4.481	4.562	3.634	4.523	18.335
satellite_50	6.934	14.146	5.139	5.179	4.226	6.255	18.375
tic-tac-toe_10	0.467	0.512	0.665	0.816	0.656	0.625	0.624
tic-tac-toe_20	0.466	0.527	0.662	0.816	0.657	0.616	0.614
tic-tac-toe_30	0.469	0.576	0.657	0.816	0.671	0.625	0.622
tic-tac-toe_50	0.470	0.606	0.661	0.816	0.685	0.647	0.618
student_10	0.258	0.247	0.897	0.897	0.293	0.293	0.514
student_20	0.258	0.247	0.536	0.288	0.289	0.294	0.514
student_30	0.295	0.285	0.537	0.287	0.292	0.303	0.515
student_50	0.298	0.268	0.536	0.307	0.289	0.322	0.515
wine_10	43.473	190.147	56.868	51.250	54.981	44.689	121.008
wine_20	46.014	171.845	46.775	46.133	51.354	53.893	106.947
wine_30	47.901	154.344	52.358	54.764	48.672	54.636	99.651
wine_50	49.698	133.361	55.972	51.696	54.857	63.823	82.517
wine-red_10	6.305	8.690	7.017	7.991	5.181	5.951	8.889
wine-red_20	7.244	9.458	8.143	7.287	7.531	7.316	9.586
wine-red_30	8.099	10.313	9.092	8.883	8.414	9.752	10.265
wine-red_50	8.381	10.114	9.566	9.636	10.265	10.408	9.879
wine-white_10	10.115	13.103	10.551	39.086	8.009	8.817	12.908
wine-white_20	10.783	14.253	11.387	42.724	9.126	10.923	13.902
wine-white_30	11.188	15.383	11.174	41.324	10.350	12.454	13.432
wine-white_50	11.214	15.699	12.028	39.693	11.935	14.184	13.373
wireless_10	4.438	6.707	4.452	5.165	3.677	3.863	7.967
wireless_20	4.735	7.082	4.554	2.325	3.892	4.491	8.464
wireless_30	4.851	8.470	4.728	1.863	4.214	4.995	8.588
wireless_50	4.473	9.132	4.677	4.918	4.496	4.629	8.470
yeast_10	0.084	0.131	0.091	0.109	0.088	0.088	0.110
yeast_20	0.083	0.126	0.089	0.082	0.087	0.091	0.099
yeast_30	0.085	0.162	0.090	0.093	0.086	0.091	0.102
yeast_50	0.088	0.156	0.096	0.088	0.096	0.092	0.102

Table 8: Root Mean Squared Error (RMSE) of the imputation performance for all considered algorithms with respect to all possible percentages of missing elements.

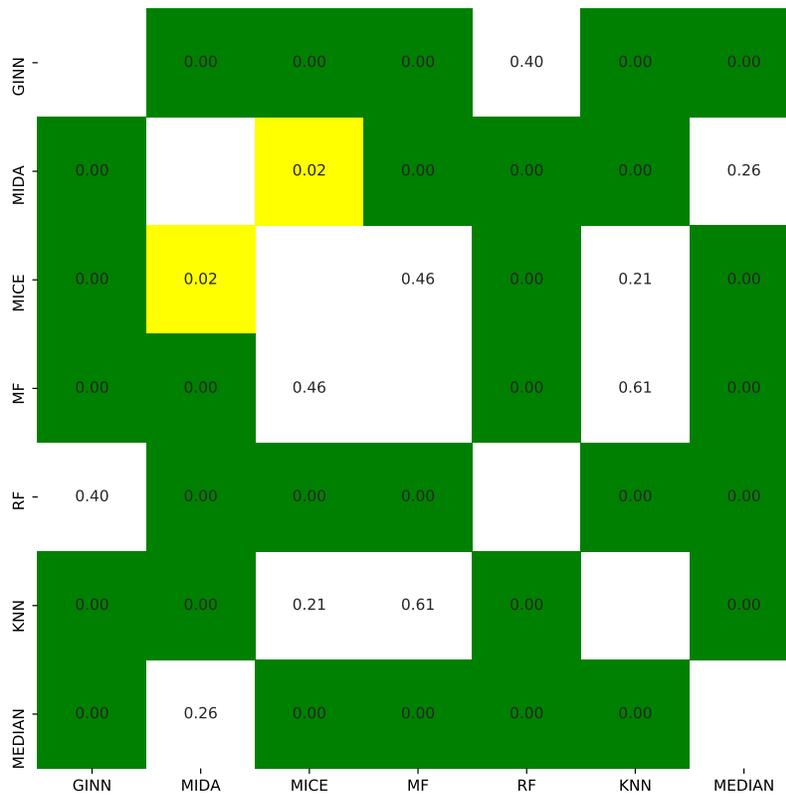
Random forest	GT	GINN	MIDA	MICE	MF	RF	KNN	MEDIAN
abalone_10	52.632	55.263	52.871	53.190	55.024	55.343	54.147	53.907
abalone_20	53.748	54.785	53.030	53.907	52.791	52.951	55.263	52.313
abalone_30	53.270	55.821	53.270	54.944	51.754	52.472	53.030	51.675
abalone_50	54.545	54.625	50.638	52.233	52.313	53.987	53.270	53.828
anuran-calls_10	93.191	93.654	92.450	91.292	90.968	92.867	92.265	91.292
anuran-calls_20	93.840	93.747	90.690	92.774	90.968	92.497	93.747	91.014
anuran-calls_30	92.728	94.442	89.810	91.848	93.191	92.635	91.061	90.412
anuran-calls_50	93.654	92.913	89.208	91.431	91.570	93.793	92.635	90.366
balance-scale_10	78.191	79.787	75.532	74.468	80.319	79.255	71.809	77.128
balance-scale_20	79.255	79.787	78.191	77.128	71.277	73.936	75.532	76.596
balance-scale_30	81.383	76.596	72.340	76.596	77.660	76.064	78.723	69.149
balance-scale_50	80.851	75.532	63.830	61.702	72.872	71.809	75.000	65.957
breast-cancer_10	97.661	95.322	97.661	95.906	95.906	98.246	95.906	95.322
breast-cancer_20	96.491	96.491	95.322	96.491	95.906	94.152	96.491	95.906
breast-cancer_30	97.076	96.491	95.906	94.737	94.152	97.076	96.491	97.661
breast-cancer_50	96.491	94.737	90.058	94.737	94.737	96.491	96.491	93.567
car-evaluation_10	71.291	71.484	71.098	70.135	69.942	69.942	70.520	70.135
car-evaluation_20	70.328	71.869	70.328	69.942	69.942	71.676	69.942	70.135
car-evaluation_30	70.520	72.832	71.484	69.942	69.942	70.713	70.520	69.942
car-evaluation_50	73.218	71.869	73.603	69.942	70.906	71.484	71.484	71.869
credit-card_10	77.933	78.453	77.887	77.873	78.160	77.947	77.880	77.893
credit-card_20	77.913	77.913	77.907	77.880	77.907	77.873	77.853	77.973
credit-card_30	77.873	77.980	77.880	77.880	77.887	77.893	77.900	77.900
credit-card_50	78.433	78.933	78.367	77.880	77.933	77.887	77.893	77.920
electrical-grid_10	97.200	97.433	95.600	99.667	97.433	88.733	99.933	98.867
electrical-grid_20	99.933	99.567	98.633	93.267	99.533	94.500	94.667	99.500
electrical-grid_30	98.533	91.100	86.633	97.567	94.900	97.867	98.700	99.367
electrical-grid_50	95.633	93.900	83.933	95.633	91.433	93.467	99.867	99.367
heart_10	81.319	84.615	82.418	80.220	83.516	78.022	73.626	82.418
heart_20	80.220	85.714	80.220	79.121	78.022	85.714	81.319	78.022
heart_30	82.418	83.516	73.626	83.516	78.022	82.418	82.418	76.923
heart_50	78.022	76.923	74.725	76.923	76.923	82.418	80.220	73.626
ionosphere_10	90.566	90.566	92.453	92.453	90.566	91.509	91.509	91.509
ionosphere_20	94.340	91.509	86.792	92.453	88.679	92.453	91.509	93.396
ionosphere_30	90.566	95.283	84.906	92.453	92.453	92.453	95.283	93.396
ionosphere_50	92.453	92.453	78.302	91.509	86.792	91.509	91.509	92.453
iris_10	88.889	95.556	93.333	93.333	91.111	93.333	93.333	95.556
iris_20	91.111	91.111	91.111	88.889	91.111	91.111	93.333	91.111
iris_30	91.111	88.889	80.000	93.333	91.111	91.111	93.333	86.667
iris_50	93.333	88.889	88.889	88.889	91.111	88.889	91.111	84.444
page-blocks_10	94.762	95.250	94.458	94.093	94.580	94.884	94.580	95.128
page-blocks_20	94.153	94.884	95.250	93.910	94.093	95.371	95.371	94.945
page-blocks_30	94.641	95.067	94.823	94.762	94.641	94.641	94.884	94.702
page-blocks_50	94.458	95.615	94.032	92.387	94.641	94.702	95.067	94.823

Random forest	GT	GINN	MIDA	MICE	MF	RF	KNN	MEDIAN
phishing_10	84.236	86.453	86.207	83.744	84.483	84.729	85.222	83.251
phishing_20	85.222	84.975	85.961	85.714	83.005	84.483	83.744	84.975
phishing_30	84.729	84.483	83.005	83.251	83.990	83.005	84.236	82.020
phishing_50	83.251	85.468	80.296	80.542	82.759	82.512	79.310	80.542
satellite_10	83.450	82.550	82.850	83.350	82.750	82.000	82.550	80.650
satellite_20	83.400	83.300	82.400	82.150	82.050	82.650	83.250	79.700
satellite_30	83.900	83.850	81.700	82.550	83.650	83.300	83.500	79.950
satellite_50	83.200	82.850	78.150	83.250	81.800	83.050	83.350	79.650
tic-tac-toe_10	74.306	75.000	70.833	67.708	68.056	67.708	69.792	68.056
tic-tac-toe_20	69.097	70.486	70.486	66.667	71.875	65.972	70.833	72.222
tic-tac-toe_30	69.444	72.569	68.056	67.708	67.361	65.972	71.181	72.917
tic-tac-toe_50	68.403	68.750	69.444	65.625	65.625	68.750	65.972	66.319
student_10	83.104	85.338	82.188	83.505	80.928	82.245	85.052	83.448
student_20	82.245	82.188	84.364	83.505	84.364	84.021	81.787	83.104
student_30	84.021	84.651	84.593	81.214	84.651	83.276	83.505	83.219
student_50	82.761	83.333	82.245	85.739	82.188	84.135	84.593	84.021
wine_10	92.593	94.444	94.444	98.148	96.296	98.148	94.444	100.000
wine_20	92.593	94.444	94.444	96.296	98.148	90.741	92.593	98.148
wine_30	96.296	98.148	92.593	98.148	94.444	98.148	94.444	94.444
wine_50	90.741	98.148	75.926	98.148	92.593	92.593	92.593	94.444
wine-red_10	61.875	63.958	59.375	60.417	61.458	63.750	62.917	61.875
wine-red_20	61.875	62.917	60.417	60.833	59.375	58.542	62.917	59.792
wine-red_30	63.333	62.708	62.083	61.667	63.333	62.292	61.458	58.125
wine-red_50	63.542	58.542	53.958	60.625	58.542	59.792	55.000	52.917
wine-white_10	51.565	54.490	52.585	53.401	51.565	53.469	53.469	53.061
wine-white_20	52.789	54.490	51.973	54.626	50.816	52.517	54.082	53.946
wine-white_30	52.109	53.673	51.361	52.517	52.313	53.265	51.020	51.837
wine-white_50	54.422	53.401	50.272	52.517	49.728	51.156	50.272	50.340
wireless_10	97.500	97.667	97.333	97.667	97.667	97.000	97.000	97.667
wireless_20	96.833	97.000	97.000	93.667	96.000	94.333	95.500	95.833
wireless_30	97.167	96.333	93.167	97.333	96.667	97.000	95.833	97.167
wireless_50	96.500	97.500	95.333	93.333	95.333	96.000	97.167	95.833
yeast_10	57.848	59.193	59.193	58.072	54.709	59.193	56.951	57.399
yeast_20	58.744	61.211	55.381	59.865	57.623	60.090	58.744	56.726
yeast_30	55.157	56.502	46.188	58.072	43.722	54.484	56.502	53.363
yeast_50	57.623	49.327	39.910	56.951	47.309	52.691	54.036	56.278

Table 9: Classification accuracy on each dataset using a random forest classifier. The model was trained over the imputed data by the analyzed algorithms for all percentage of missing elements.

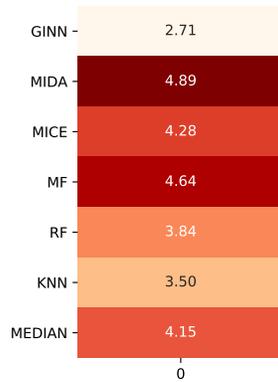


(a) Rankings

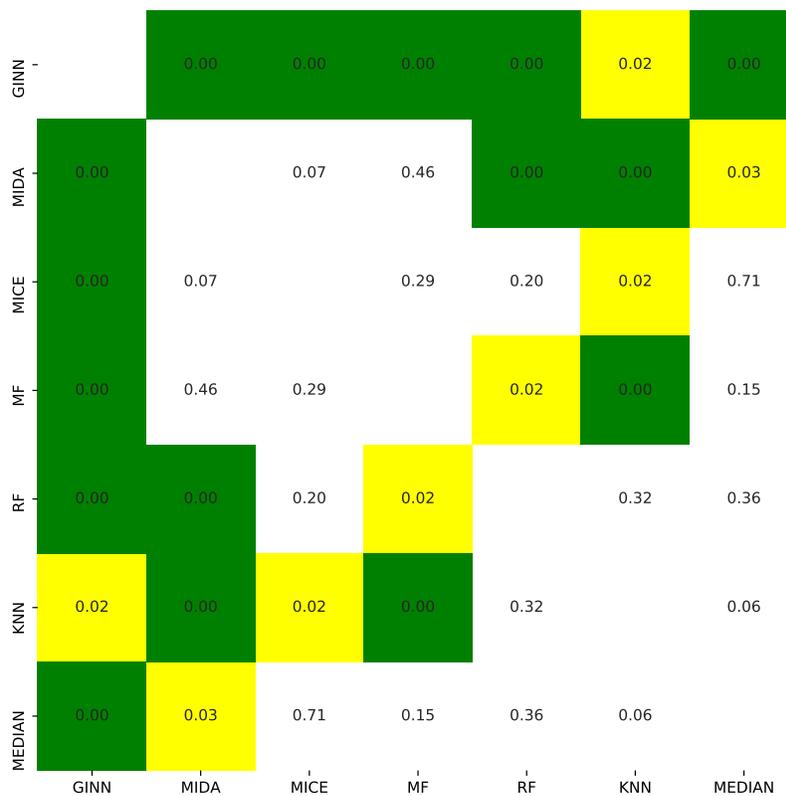


(b) P-values

Figure 7: (a) Rankings of the algorithms in Section 5.1 with respect to RMSE, averaged over all datasets. (b) Corresponding p-values of a set of post-hoc Nemenyi tests. Green shows significant differences at 0.01, yellow significant differences at 0.05.

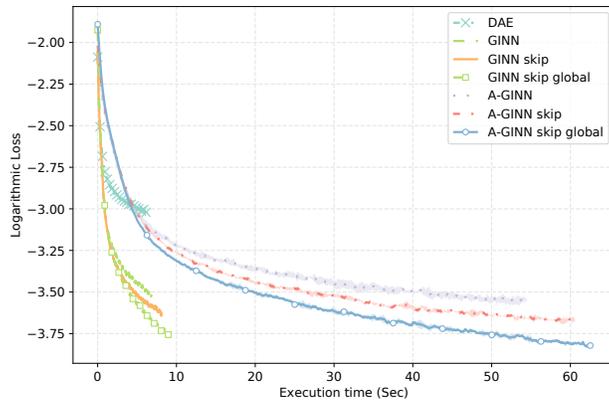


(a) Rankings

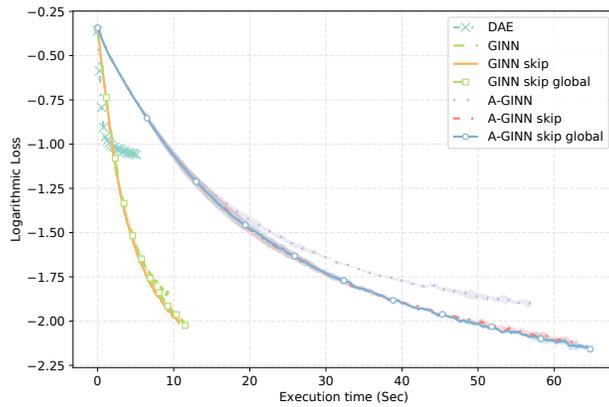


(b) P-values

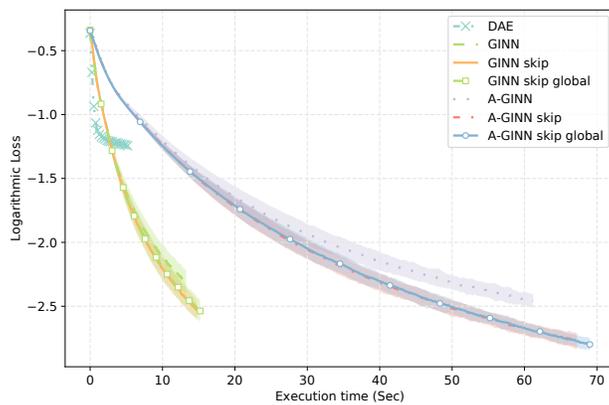
Figure 8: (a) Rankings of the algorithms in Section 5.2 with respect to the accuracy of random forest, averaged over all datasets. (b) Corresponding p-values of a set of post-hoc Nemenyi tests. Green shows significant differences at 0.01, yellow significant differences at 0.05.



(a)

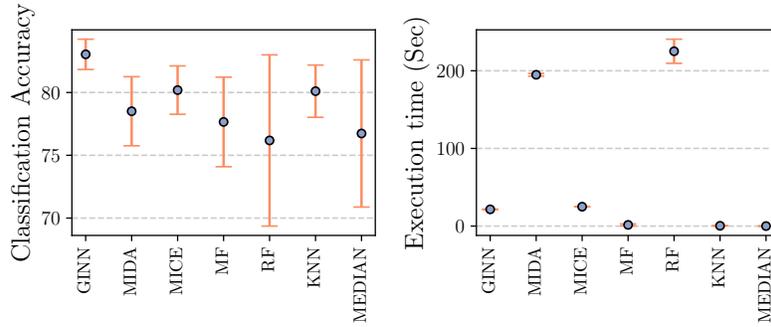


(b)

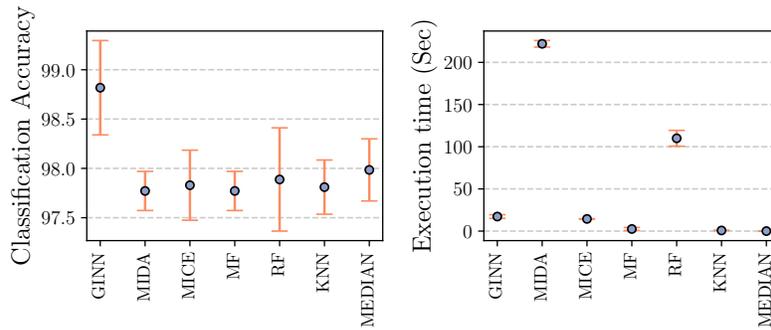


(c)

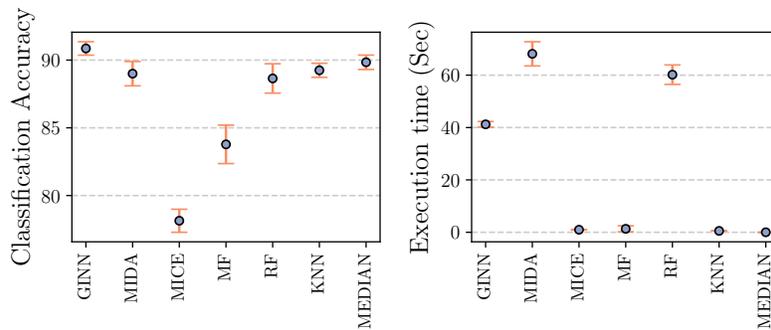
Figure 9: Convergence of the logarithmic loss function with respect to a fixed computational budget on the Ionosphere (a), Tic-Tac-Toe (b) and Phishing (c) datasets for each variant described in the ablation study of Section 5.3.



(a) Mammographic mass

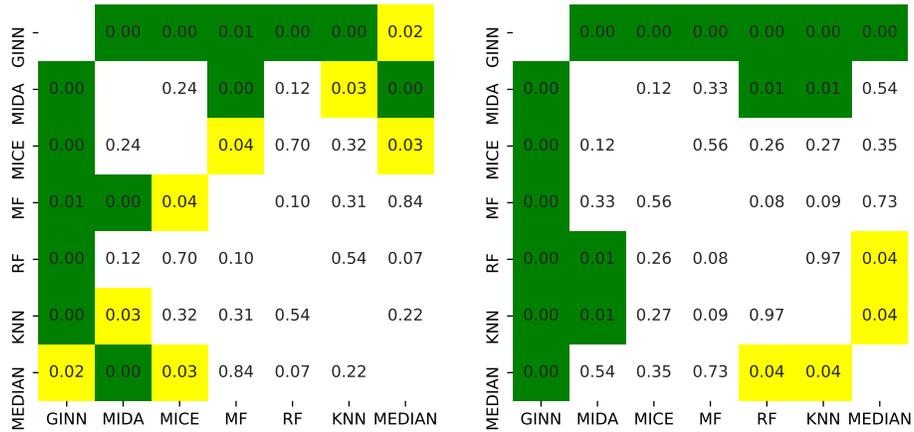


(b) Cervical cancer



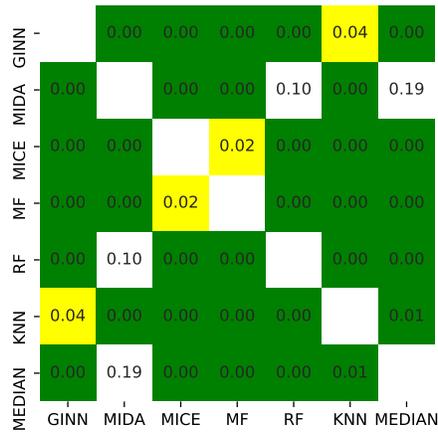
(c) Air quality

Figure 10: Random forest classification accuracy and imputation times on (a) Mammographic mass (b) Cervical cancer and (c) Air quality datasets.



(a) Cervical cancer

(b) Mammography mass



(c) Air quality

Figure 11: Corresponding p-values of a set of post-hoc Nemenyi tests on the rankings of Table 6 in the paper relative to the downstream classification/regression tasks on the three datasets with pre-existing missing data: (a) Cervical cancer (13% of missing elements), (b) Mammographic mass (4% of missing elements) and (c) Air quality (13% of missing elements). Green shows significant differences at 0.01, yellow significant differences at 0.05.