

Siamese Autoencoder-based Approach for Missing Data Imputation

Ricardo Cardoso Pereira¹, Pedro Henriques Abreu¹, and Pedro Pereira Rodrigues²

¹ Centre for Informatics and Systems of the University of Coimbra, Department of Informatics Engineering, University of Coimbra, 3030-290 Coimbra, Portugal

`{rdpereira,pha}@dei.uc.pt`

² Center for Health Technology and Services Research, Faculty of Medicine (MEDCIDS), University of Porto, 4200-319 Porto, Portugal

`pprodrigues@med.up.pt`

Abstract. Missing data is an issue that can negatively impact any task performed with the available data and it is often found in real-world domains such as healthcare. One of the most common strategies to address this issue is to perform imputation, where the missing values are replaced by estimates. Several approaches based on statistics and machine learning techniques have been proposed for this purpose, including deep learning architectures such as generative adversarial networks and autoencoders. In this work, we propose a novel siamese neural network suitable for missing data imputation, which we call Siamese Autoencoder-based Approach for Imputation (SAEI). Besides having a deep autoencoder architecture, SAEI also has a custom loss function and triplet mining strategy that are tailored for the missing data issue. The proposed SAEI approach is compared to seven state-of-the-art imputation methods in an experimental setup that comprises 14 heterogeneous datasets of the healthcare domain injected with Missing Not At Random values at a rate between 10% and 60%. The results show that SAEI significantly outperforms all the remaining imputation methods for all experimented settings, achieving an average improvement of 35%.

Keywords: Missing Data · Imputation · Siamese Autoencoder · Missing Not At Random

1 Introduction

Missing data can be described by the absence of values in the instances of a dataset. It is one of the most common data issues, affecting most real-world domains. The existence of missing values has a deep impact on the conclusions that can be drawn from the data [10]. Within machine learning, the majority of the classification and regression models cannot cope with missing data, or suffer a decrease in their performance. In domains such as healthcare, missing values can compromise the results to the point where the study becomes unfeasible [12]. However, different types of missing values may have different impacts. Missing

data can be classified according to three different mechanisms which are related to the missingness causes [15, 8]:

- Missing Completely At Random (MCAR), which is described by Eq. 1, where $P(x)$ is the probability function, R is the binary matrix of the missing values in Y , Y is a combination of the observed values Y_{obs} and the missing values Y_{mis} , and ψ are the parameters of the missing data model. In this mechanism the missingness causes are purely random, and therefore only depended on ψ . An example would be if someone simply forgot to fill in a field in a questionnaire.

$$P(R = 0|Y_{obs}, Y_{mis}, \psi) = P(R = 0|\psi) \quad (1)$$

- Missing At Random (MAR), which is described by Eq. 2 and the missingness nature is related to observed data in one or more features, therefore depending on both Y_{obs} and ψ . For example, people may not have results for a specific medical exam because they are too young to be doing such an exam.

$$P(R = 0|Y_{obs}, Y_{mis}, \psi) = P(R = 0|Y_{obs}, \psi) \quad (2)$$

- Missing Not At Random (MNAR), where the missingness causes are unknown since the missing values are related with themselves or to external data not available. As a consequence, the missingness depends on Y_{obs} , Y_{mis} and ψ . For example, people that drink a lot of alcohol may be apprehensive about answering how many drinks they have per day.

A solution often used to address the missing data issue is imputation, which consists in the replacing the missing values by estimates [19]. There are several statistical and machine learning-based methods to perform imputation, and each one may be more suitable for a specific missing mechanism [10]. However, most approaches are tailored for MCAR and MAR, while MNAR still is the less addressed mechanism with few solutions, mostly because of its relation with unknown data. The proposal of new approaches that perform well under the MNAR mechanism is an important open challenge when considering that most real-world contexts suffer from this mechanism, including critical domains such as healthcare [13].

Recently, several deep learning architectures have been explored to perform imputation assuming the different mechanisms. The state-of-the-art architectures in this scope are generative adversarial networks (GAN) [21], denoising autoencoders (DAE) [20, 5] and variational autoencoders (VAE) [11, 13]. In this work, we explore the use of siamese networks [6] to perform missing data imputation. We propose a new model called Siamese Autoencoder-based Approach for Imputation (SAEI), which extends and adapts the vanilla siamese network for the imputation task by adding a deep autoencoder architecture, a custom loss function and a custom triplet mining strategy. To the best of our knowledge, this is the first time siamese networks have been used for this purpose. We compared our SAEI model with a baseline of seven state-of-the-art imputation methods in a experimental setup encompassing 14 datasets of the healthcare domain injected with MNAR values at different missingness levels (10% to 60% missing

rates). The achieved results proved that our SAEI model significantly outperformed the remaining baseline methods in all experimented settings, achieving an average improvement of 35%.

The remainder of the paper is organized in the following way: Section 2 presents related work on the missing data field; Section 3 describes with detail the proposed SAEI model; Section 4 presents the experimental setup used to validate the effectiveness of our SAEI model; Section 5 analyzes and discusses the results obtained from the experiments; and Section 6 presents our conclusions and the future work.

2 Related Work

As previously stated, the missing data problem cannot be neglected, otherwise it will have a major impact in any tasks performed with the data. An approach often used to deal with this issue is to remove the missing values, which is called case deletion. Two different deletion approaches can be applied [10]: listwise deletion, where all instances containing at least one missing value are eliminated, which is a very simple procedure but can lead to the loss of a considerable amount of information, being for that reason only recommended in big data scenarios or when the instances with missing values are less than 5% of the data; and pairwise deletion, where only the instances containing missing values for the features of interest are deleted, which suffers from the same problems of listwise deletion but reduces the loss of information. To apply any deletion strategies, the missing mechanism must be considered since deleting data may remove relations between instances and features, and for that reason this approach should only be applied with missing values under MCAR [10].

When deleting instances is not a suitable option, imputation tends to be the preferred strategy [19]. The key idea is to generate plausible new estimates to replace the missing values. Among the existent methods, the statistical-based ones are frequently used [8]. A common approach is to use the mean or mode of the feature containing missing values for numeric and nominal variables, respectively [10]. Another strategy is to use a regression to model one or more dependent variables (the ones containing missing values) using as independent variables the remaining features of the dataset [10]. The fitting process must only consider the instances that are complete for the independent features, and the type of regression (e.g., linear or non-linear) must be chosen taking in consideration the nature of the data. Both latter approaches are said to be single imputation strategies, since only one value is used for the imputation of each missing value. For this reason, the imputation results may be biased since the uncertainty of the generated values is not accounted for. To address this issue multiple imputation strategies can be applied. The key idea is to perform the imputation M times, generating different but complete results each time (different imputation methods may be used). The M complete datasets are then analyzed and combined, being the result of this combination the final dataset [16]. The state-of-the-art method that uses this approach is the Multiple Imputation by

Chained Equations (MICE) [4]. It creates a series of regressions where each one is modeled for a different variable with missing values, meaning that each feature is modeled conditionally upon the other features. The process is repeated throughout several iterations until the parameters of the regressions converge to stable values [1].

Other approaches based on statistical and algebra concepts often used are matrix completion methods, such as the SoftImpute. The goal is to perform the imputation by finding a low-rank matrix that is an approximation of the original one. The process usually relies on matrix factorization, where a matrix with (m, n) dimensions is decomposed into two matrices with dimensions (m, k) and (k, n) , where k is the rank and must be smaller than m and n . The idea is to find the latent features that better describe the available values. The low-rank approximation matrix can finally be obtained by multiplying the resulting matrices from the decomposition process [18].

The imputation task can also be performed by using machine learning models. In theory, any algorithm that can be trained to predict new values is suitable for imputation [8]. Generating new estimates for numerical features may be seen as a regression problem, or a classification problem for nominal features.

An algorithm often used for this purpose is the k -nearest neighbors (KNN) [2, 7]. It tries to find the k most similar instances to the one that contains missing values using only the features that are complete. To calculate this similarity a distance function is used, and it must be chosen taking in consideration the data type: the euclidean distance is suitable for numeric data but not for nominal. For this latter case the data can be transformed through one-hot encoding or a different distance function must be used (e.g., hamming distance works with nominal data) [3]. When $k > 1$ the values of the neighbors' instances must be combined to produce the new value. For numerical data a common approach used is the simple or weighted mean, and for categorical values a vote of majority may be applied.

Neural networks are also often used for imputation, particularly autoencoders and generative adversarial networks (GAN). Autoencoders are a type of neural network that learns a representation of the data from the input layer and tries to reproduce it at the output layer. Among its variants, the denoising autoencoder (DAE) is the one more often used for missing data imputation because it is designed to recover noisy data, which can exist due to data corruption via some additive mechanism or by the introduction of missing values [5]. However, the variational autoencoder (VAE) has recently been used for the same purpose. This architecture learns the multidimensional parameters of a Gaussian distribution and generates the estimates to replace the missing values by sampling from that distribution [13]. Regarding GANs, these are trained through an adversarial process and are based on a generative model that learns the data distribution and a discriminative model that outputs the probability of a sample being from the training data rather than the generative model [9]. GANs are directly applied to the missing data imputation task through the generative adversarial

imputation nets (GAIN) [21], where the generator performs the imputation and the discriminator tries to distinguish between the original and the imputed data.

A type of neural network that has not yet been used for missing data imputation is the siamese network. This architecture is based on two or more sub-networks that share the same architecture and parameters. The goal is to output embedding representations in a way that similar concepts would have a similar latent space embedding [6, 17]. In this work, we propose an extension to the siamese network tailored for missing data imputation, which is, to the best of our knowledge, the first time that such architecture is being used for this purpose. Since this network uses a triplet loss function, which is based on distances between anchor, positive and negative samples, it can learn from reduced quantities of data (especially when compared to other neural network-based approaches). This quality makes this type of network feasible to be used with both lower and higher missing rates. Furthermore, the number of complete samples (i.e., without containing missing values) needed to train the model can be rather small when compared to other network-based imputation models. We leverage these characteristics of the siamese networks and propose an adapted method for missing data imputation with a novel deep autoencoder architecture, custom loss function and custom triplet mining strategy.

3 Siamese Autoencoder-based Approach for Imputation

The Siamese Autoencoder-based Approach for Imputation (SAEI) is an extension of a vanilla siamese network, and it is tailored for missing data imputation by comprising three adaptations:

- A deep autoencoder architecture that allows the network to reproduce the input data at the output layer in an unsupervised fashion;
- A custom loss function that includes both the distance-based triplet loss and the reconstruction error of the autoencoder component of the network;
- A custom triplet mining strategy that was designed specifically for the missing data issue by creating hard triplets based on the existing missing values.

These three adaptations are independently described in the next subsections.

3.1 Deep Autoencoder Architecture

The architecture of our SAEI model is roughly inspired in the well-know ZFNet [22], although it presents several changes motivated by it being aimed at tabular data. Fig. 1 and Fig. 2 depict graphical descriptions of the encoder and decoder networks, respectively. The encoder network is composed of two one-dimensional convolutional layers with 16 filters, ReLU as the activation function, and kernel sizes of five and three. Such layers are followed by max-pooling layers with two strides, which are then followed by two regularization layers that perform batch normalization and dropout at a rate of 25%. Moreover, a residual connection is also used to skip the second convolutional layer. Finally, the output

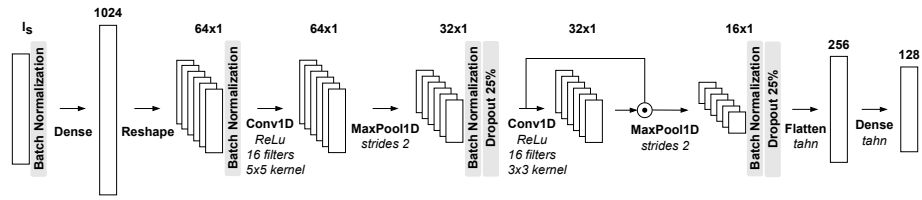


Fig. 1. SAEI encoder architecture. I_s represents the input shape. The encoder network is composed of two one-dimensional convolutional layers with 16 filters followed by max-pooling layers with two strides. Regularization is performed through batch normalization and dropout at a rate of 25%. A residual connection is also used to skip the second convolutional layer. The latent output is obtained from a dense layer with 128 units.

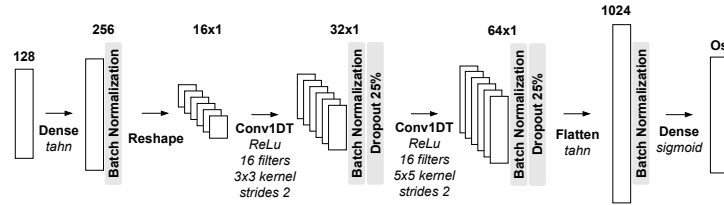


Fig. 2. SAEI decoder architecture. O_s represents the output shape. The decoder network is symmetric to the encoder. Therefore, it has the same layers and regularization but in reverse order. The residual connection is not applied by the decoder.

from the last pooling layer is flattened and passed through a hyperbolic tangent activation function, and the latent output is obtained from a final dense layer with 128 units and the latter activation function. The decoder network is symmetric to the encoder, presenting the same architecture but in reverse order (without the residual connection). To perform the deconvolution operation, the one-dimensional convolutional and the max pooling layers are replaced by one-dimensional transposed convolutional layers with two strides. The output dense layer uses the sigmoid activation function so that the data can be normalized within $[0, 1]$.

Convolutional layers operate based on the spatial positioning of the data. In other words, the position where each value is placed is relevant for the feature extraction process. Tabular data does not present this behavior because the features' positions are irrelevant. To surpass this limitation, the SAEI model feeds the input data to a dense layer with 1024 units and without activation function. The purpose of this layer is to learn an abstract representation of the input data where the spatial relation between the new abstract features is meaningful. Therefore, our SAEI model delegates the task of learning the spatial structure of the features to the network. The first convolutional layer is then fed with the output of the mentioned dense layer.

3.2 Custom Loss Function

One of the original loss functions proposed to train a siamese network was the triplet loss. It tries to minimize the distance between an anchor and a positive sample which represent the same concept, while maximizing the distance between that same anchor and a negative sample of a different concept. The formulation is presented in Eq. 3, where $f(x)$ is a function of the embedding representation, N is the number of triples composed by an anchor (x_i^a), a positive sample (x_i^p) and a negative one (x_i^n) [17]. Furthermore, α is a margin used to ensure a minimum distance between positive and negative samples.

$$TL_i = \sum_i^N \left[\|f(x_i^a) - f(x_i^p)\|_2^2 - \|f(x_i^a) - f(x_i^n)\|_2^2 + \alpha \right] \quad (3)$$

Using the triplet loss function is ideal for a vanilla siamese network since the goal is for the model to distinguish between positive and negative concepts in comparison to the anchor, while outputting embedding representations that incorporate such differences. However, our SAEI model relies on an autoencoder architecture that outputs a reconstruction of the input data based on the anchor latent representation. Therefore, this reconstruction component must also be reflected in the loss function of the model. Eq. 4 presents our custom loss function, where the triplet loss (TL_i) from Eq. 3 is added to the mean squared error between the anchor and positive sample, similarly to what would happen in a denoising autoencoder (i.e., the anchor is the corrupted version of the ground truth represented by the positive sample).

$$TL_i + \sum_i^N (x_i^{\hat{a}} - x_i^p)^2 \quad (4)$$

3.3 Custom Triplet Mining

The triplet selection is a key step for successfully training a siamese network. To ensure that the network is able to learn how to distinguish between the positive and negative concepts, it is imperative to select triplets that the network is unable to differentiate before being trained. These triplets are composed by the so-called hard positives and hard negatives, which are samples that violate the constraint of the triplet loss, presented in Eq. 5 [17]. If the network is trained with easy triples where the triplet loss constraint is already satisfied before training, it would not gain the capacity of distinguishing between positive and negative concepts.

$$\|f(x_i^a) - f(x_i^p)\|_2^2 + \alpha < \|f(x_i^a) - f(x_i^n)\|_2^2 \quad (5)$$

When extending this type of network and its respective training procedure to address missing data imputation, the definitions of anchor, negative and positive samples must be redefined within the missing data scope. Our SAEI model proposes the following definitions:

- Anchors are the samples containing missing values, which are pre-imputed with the mean of the available values in each feature (or the mode in categorical features). This pre-imputation step is required since neural networks are unable to be trained with data containing missing values.
- Positives are the original samples without containing missing values (i.e., the ground truth of the anchor samples). This implies that we must have a portion of complete data to train the SAEI model, but such assumption is common to all deep learning-based imputation methods. Furthermore, the anchor samples are created by artificially generating missing values in the positive samples, according to a pre-established missing data mechanism.
- Negatives are the same as anchor samples but with the missing values being replaced by Gaussian noise sampled according to Eq. 6, where D represents the entire dataset and λ represents the variance of the noise domain.

$$\mathcal{N}\left(\frac{\max(D) - \min(D)}{2}, \lambda\right) \quad (6)$$

The rationale behind the mean value is to use the midpoint within the domain of the data, therefore keeping the Gaussian noise centered on that domain. Such strategy works specially well when the data is normalized within a specific range (e.g., $[0, 1]$). Moreover, the λ parameter acts as a control variable to define how hard or easy should the triplets be: a low λ value leads to Gaussian noise mostly or completely contained within the data domain, therefore leading to harder triplets since the negative sample is likely to be partially overlapped with the positive one; on the other hand, a high λ value increases the noise domain and creates a negative sample that is more dissimilar to the positive, therefore creating easier triplets where the negative and positive samples are immediately distinguishable by the model. As a consequence, this parameter should be defined aiming to generate hard triplets while considering for the data domain.

4 Experimental Setup

The quality of the imputation results achieved by the SAEI model was assessed by comparing it with a baseline of seven state-of-the-art imputation methods introduced in Section 2: kNN with $k = 5$, Mean/Mode, MICE with 100 iterations, SoftImpute with 100 iterations, DAE, VAE and GAIN. The DAE and VAE were defined according to the following architecture and hyperparameters: a hidden layer with half of the input dimension and ReLU as the activation function (although VAE has two additional layers for the Gaussian distribution parameters), Adam as the optimizer with a learning rate of 0.001, Sigmoid as the activation function of the output layer (forcing the data to be normalized within $[0, 1]$), batches of 64 samples, 200 training epochs, a dropout rate of 10% for regularization, Mean Squared Error as the reconstruction loss, and early stop and learning rate reduction by 80% if the validation loss does not improve over 100 epochs. The VAE loss function also includes the Kullback–Leibler divergence for

regularization. Furthermore, both DAE and VAE require pre-imputation, which is performed with the mean/mode of the features. The GAIN method was used with the hyperparameters proposed by its authors. Our SAEI model follows the architecture described in Section 3.1 while all the remaining hyperparameters are the same used by the DAE and VAE. Moreover, the Gaussian noise of the negative instances was sampled from $\mathcal{N}(0.5, 0.05^{0.5})$, and the margin of the triplet loss function was set to $\alpha = 0.2$ as the original authors proposed. With the exception of GAIN, the described hyperparameters of all methods were defined through a grid search process. This is a standard procedure that aims to achieve hyperparameters that conform to common use cases. Moreover, the DAE and VAE were implemented with the Keras library, the GAIN implementation was obtained from its original authors³, and the remaining methods were obtained and used from the Scikit-learn library. Our SAEI model was also coded with the Keras library and is available on GitHub⁴.

The experiments were conducted over 14 datasets of the healthcare domain which cover different types of clinical data that was collected for several pathologies. We chose to cover this medical domain since it often suffers from the missing data issue, particularly missing values under the MNAR mechanism, which creates deep challenges to any subsequent analysis or task performed with the data [12]. In this health context, each instance usually represents a group of values collected for a patient, and each value belongs to a feature that could be a measurement, an exam result, or any other medical input. The missing values may appear at any feature and/or instance. All the 14 datasets are public and available at the UC Irvine Machine Learning⁵ and Kaggle⁶ repositories, and they present heterogeneous characteristics as seen in Table 1.

The datasets were normalized within $[0, 1]$ and split into train and test sets with 70% and 30% of the instances, respectively. The scaler learns the minimum and maximum values from the train set and transforms both sets. This strategy ensures the test data is not biased with information from the training data. However, in the presence of high missing rates (usually above 50%), the test set may contain values unseen by the scaler, which leads to the normalization boundaries being slightly extended from the expected $[0, 1]$ domain. Also, the autoencoder-based methods use 20% of the train set as the validation set. Furthermore, the categorical nominal features were transformed through the one-hot encoding procedure so that they can be supported by all imputation methods. Finally, in order to be able to calculate the imputation error of the experimented methods, the datasets must be complete and the test set must be injected with artificially generated missing values so that we can compare the estimated values with the original ones (i.e., ground truth). We choose to generate MNAR values since it is the hardest mechanism to address and the one more often found in real-world contexts such as healthcare [12, 13]. Our generation strategy is based

³ <https://github.com/jsyoon0823/GAIN>

⁴ <https://github.com/ricardodcperreira/SAEI>

⁵ <https://archive.ics.uci.edu/ml>

⁶ <https://www.kaggle.com/datasets>

Table 1. Datasets characteristics.

Dataset	# Instances	# Features	
		Categorical	Continuous
diabetic-retinopathy	1151	4	16
ecoli	336	1	7
ctg-2c	2126	1	21
new-thyroid-N-vs-HH	215	1	5
kala-azar	68	2	5
immunotherapy	90	3	5
saheart	462	2	8
bc-coimbra	116	1	9
cleveland-0-vs-4	173	1	13
newthyroid-v1	185	1	5
biomed	194	1	5
cryotherapy	90	3	4
thyroid-3-vs-2	703	1	21
pima	768	1	8

on removing the smaller values of several features at once (upon a certain missing rate) in a multivariate procedure. Therefore, the missing rate is defined for the entire dataset and the imputation is performed on all features simultaneously. The obtained results were evaluated with the Mean Absolute Error (MAE) between the estimates and the ground truth values. Regarding the missing rate, we considered four levels of missingness: 10%, 20%, 40% and 60%.

To avoid the impact of stochastic behaviors in the results, the experiment was executed 30 independent times, with the train/test split being randomly performed in every run. The MAE results of the experiment are the average of the 30 runs. Moreover, each run was executed in a computer with the following specifications: Windows 11, CPU AMD Ryzen 5600X, 16GB RAM, and GPU NVIDIA GeForce GTX 1060 6GB. The time complexity of each imputation method was not measured, but all methods were computed in a feasible amount of time.

5 Results

The imputation results obtained from the experimental setup are displayed with detail in Table 2. The MAE results (average and standard deviation of the 30 independent runs) are individually displayed for each dataset and missing rate. Furthermore, the overall results for each imputation method and per missing rate (i.e., considering all datasets) are displayed in Fig. 3.

In an overall analysis, our SAEI model clearly outperforms all the remaining imputation methods, achieving smaller MAE values for every dataset and missing rate considered, as seen in Table 2. In fact, the average improvement of the SAEI model in comparison to the remaining baseline methods is 35%, peaking at the 20% missing rate with an improvement of 42%. For the lower

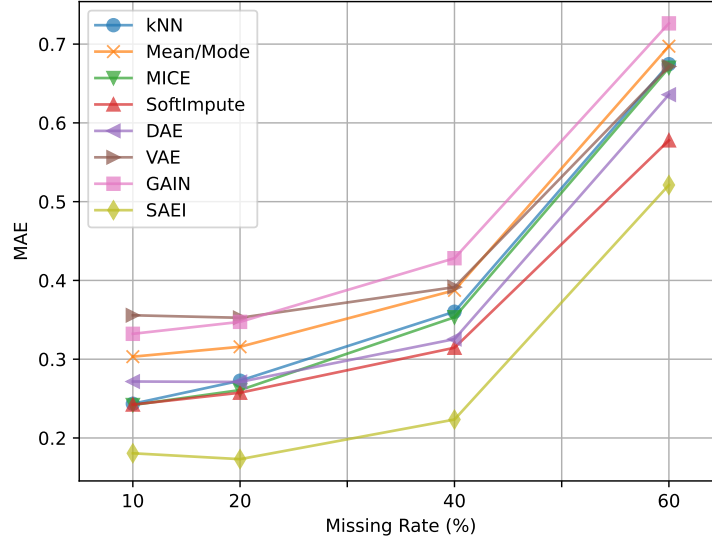


Fig. 3. Overall Mean Absolute Error of all imputation methods per missing rate. Our SAEI model significantly outperforms the remaining methods in all settings.

were not met, so the data was transformed into rankings using the Ordered Quantile normalization [14]. Such assumptions were also ensured for the data subgroups. The obtained p -values show that the results are statistically significant for all factors, with $p < 0.001$. Additionally, to validate if our SAEI model outperformed the remaining methods with a statistical significance of 5%, the post-hoc Tukey’s HSD test was applied to this factor. The obtained p -values show that the SAEI model significantly outperformed all the baseline of imputation methods, with $p < 0.001$ in all scenarios.

6 Conclusions

In this work we propose a new model for missing data imputation called Siamese Autoencoder-based Approach for Imputation (SAEI), which is an extension of the vanilla siamese networks adapted for this imputation task. The model incorporates three main adaptations: a deep autoencoder architecture that is tailored for missing data reconstruction, a custom loss function that encompasses both the triplet and reconstruction losses, and a triplet mining strategy tailored for the missing data issue that is capable of generating hard triplets that are meaningful for the training procedure. To the best of our knowledge, this is the first time a siamese architecture is being used and adapted for missing data imputation. We

compared our SAEI model with seven state-of-the-art imputation methods from both statistical and machine learning backgrounds, in an experimental setup that used 14 datasets of the healthcare domain injected with MNAR values in a multivariate fashion at 10%, 20%, 40% and 60% missing rates. Our SAEI model significantly outperformed all the baseline methods used for comparison in all the datasets and missing rates, with a statistical significance of 5%, achieving an average improvement of 35%.

In the future we want to extend this work to test other missing data mechanisms, and we want to explore automatic approaches to evolve our deep autoencoder architecture so it can be even further optimized.

Acknowledgements This work is supported in part by the FCT - Foundation for Science and Technology, I.P., Research Grant SFRH/BD/149018/2019. This work is also funded by the FCT - Foundation for Science and Technology, I.P./MCTES through national funds (PIDDAC), within the scope of CISUC R&D Unit - UIDB/00326/2020 or project code UIDP/00326/2020.

References

1. Azur, M.J., Stuart, E.A., Frangakis, C., Leaf, P.J.: Multiple imputation by chained equations: what is it and how does it work? *International journal of methods in psychiatric research* **20**(1), 40–49 (2011)
2. Batista, G., Monard, M.: Experimental comparison of k-nearest neighbor and mean or mode imputation methods with the internal strategies used by c4. 5 and cn2 to treat missing data. *University of Sao Paulo* **34** (2003)
3. Batista, G.E., Monard, M.C., et al.: A study of k-nearest neighbour as an imputation method. *HIS* **87**(251-260), 48 (2002)
4. Buuren, S.v., Groothuis-Oudshoorn, K.: mice: Multivariate imputation by chained equations in r. *Journal of Statistical Software* pp. 1–68 (2010)
5. Charte, D., Charte, F., García, S., del Jesus, M.J., Herrera, F.: A practical tutorial on autoencoders for nonlinear feature fusion: Taxonomy, models, software and guidelines. *Information Fusion* **44**, 78–96 (2018)
6. Chicco, D.: Siamese neural networks: An overview. *Artificial Neural Networks* pp. 73–94 (2021)
7. García-Laencina, P.J., Abreu, P.H., Abreu, M.H., Afonso, N.: Missing data imputation on the 5-year survival prediction of breast cancer patients with unknown discrete values. *Computers in biology and medicine* **59**, 125–133 (2015)
8. García-Laencina, P.J., Sancho-Gómez, J.L., Figueiras-Vidal, A.R.: Pattern classification with missing data: a review. *Neural Computing and Applications* **19**(2), 263–282 (2010)
9. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. *Communications of the ACM* **63**(11), 139–144 (2020)
10. Little, R.J., Rubin, D.B.: *Statistical Analysis with Missing Data*, vol. 793. John Wiley & Sons (2019)
11. McCoy, J.T., Kroon, S., Auret, L.: Variational autoencoders for missing data imputation with application to a simulated milling circuit. *IFAC-PapersOnLine* **51**(21), 141–146 (2018)

12. Peek, N., Rodrigues, P.P.: Three controversies in health data science. *International Journal of Data Science and Analytics* **6**(3), 261–269 (2018)
13. Pereira, R.C., Abreu, P.H., Rodrigues, P.P.: Partial multiple imputation with variational autoencoders: Tackling not at randomness in healthcare data. *IEEE Journal of Biomedical and Health Informatics* **26**(8), 4218–4227 (2022)
14. Peterson, R.A., Cavanaugh, J.E.: Ordered quantile normalization: a semiparametric transformation built for the cross-validation era. *Journal of Applied Statistics* pp. 1–16 (2019)
15. Rubin, D.B.: Inference and missing data. *Biometrika* **63**(3), 581–592 (1976)
16. Rubin, D.B.: *Multiple Imputation for Nonresponse in Surveys*, vol. 81. John Wiley & Sons (2004)
17. Schroff, F., Kalenichenko, D., Philbin, J.: Facenet: A unified embedding for face recognition and clustering. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 815–823 (2015)
18. Udell, M., Horn, C., Zadeh, R., Boyd, S., et al.: Generalized low rank models. *Foundations and Trends in Machine Learning* **9**(1), 1–118 (2016)
19. Van Buuren, S.: *Flexible imputation of missing data*. Chapman and Hall/CRC (2018)
20. Vincent, P., Larochelle, H., Bengio, Y., Manzagol, P.A.: Extracting and composing robust features with denoising autoencoders. In: *Proceedings of the 25th International Conference on Machine learning*. pp. 1096–1103 (2008)
21. Yoon, J., Jordon, J., Schaar, M.: Gain: Missing data imputation using generative adversarial nets. In: *International Conference on Machine Learning*. pp. 5689–5698. PMLR (2018)
22. Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional networks. In: *European conference on computer vision*. pp. 818–833. Springer (2014)