

Handling missing values in healthcare data: A systematic review of deep learning-based imputation techniques

Mingxuan Liu^{1#}, Siqu Li^{1#}, Han Yuan¹, Marcus Eng Hock Ong^{2,3}, Yilin Ning¹, Feng Xie^{1,2}, Seyed Ehsan Saffari^{1,2}, Victor Volovici⁴, Bibhas Chakraborty^{1,2,5,6}, Nan Liu^{1,2,7,8*}

¹ Centre for Quantitative Medicine, Duke-NUS Medical School, Singapore

² Programme in Health Services and Systems Research, Duke-NUS Medical School, Singapore

³ Department of Emergency Medicine, Singapore General Hospital, Singapore

⁴ Department of Neurosurgery, Erasmus MC University Medical Center, Rotterdam, The Netherlands

⁵ Department of Statistics and Data Science, National University of Singapore, Singapore

⁶ Department of Biostatistics and Bioinformatics, Duke University, Durham, NC, USA

⁷ SingHealth AI Office, Singapore Health Services, Singapore

⁸ Institute of Data Science, National University of Singapore, Singapore

These authors contributed equally

* Corresponding author: Nan Liu, Centre for Quantitative Medicine, Duke-NUS Medical School, 8 College Road, Singapore 169857, Singapore. Phone: +65 6601 6503. Email: liu.nan@duke-nus.edu.sg

Abstract

Objective

The proper handling of missing values is critical to delivering reliable estimates and decisions, especially in high-stakes fields such as clinical research. The increasing diversity and complexity

of data have led many researchers to develop deep learning (DL)-based imputation techniques. We conducted a systematic review to evaluate the use of these techniques, with a particular focus on data types, aiming to assist healthcare researchers from various disciplines in dealing with missing values.

Materials and Methods

We searched five databases (MEDLINE, Web of Science, Embase, CINAHL, and Scopus) for articles published prior to August 2021 that applied DL-based models to imputation. We assessed selected publications from four perspectives: health data types, model backbone (i.e., main architecture), imputation strategies, and comparison with non-DL-based methods. Based on data types, we created an evidence map to illustrate the adoption of DL models.

Results

We included 64 articles, of which tabular static (26.6%, 17/64) and temporal data (37.5%, 24/64) were the most frequently investigated. We found that model backbone(s) differed among data types as well as the imputation strategy. The "integrated" strategy, that is, the imputation task being solved concurrently with downstream tasks, was popular for tabular temporal (50%, 12/24) and multi-modal data (71.4%, 5/7), but limited for other data types. Moreover, DL-based imputation methods yielded better imputation accuracy in most studies, compared with non-DL-based methods.

Conclusion

DL-based imputation models can be customized based on data type, addressing the corresponding missing patterns, and its associated "integrated" strategy can enhance the efficacy of imputation, especially in scenarios where data is complex. Future research may focus on the portability and fairness of DL-based models for healthcare data imputation.

Keywords: Missing value, Imputation, Deep learning, Neural networks, Healthcare

Abbreviations: DL, Deep Learning; AE, Auto Encoder; DAE, Deep Auto Encoder; GAN, Generative Adversarial Network; LSTM, Long Short-term Memory; MLP, Multilayer Perceptron; RNN, Recurrent Neural Networks

1. Introduction

Healthcare data has emerged in diverse formats in the era of big data. Personalized health monitoring devices, for instance, enable the collection of data tailored to an individual's activities on a daily basis. Likewise, rapidly evolving laboratory techniques generate vast amounts of sequencing data. However, these new data formats are more susceptible to the problem of missing values than traditional tabular clinical data collected in a prospective observational or randomized trial.

Missing values cast a shadow on data analysis: they can reduce prediction power and lead to bias in downstream decision-making[1, 2], which is particularly problematic in high-fidelity decision-making situations, such as those in healthcare. Simple complete data analysis or simple imputation (mean, median, or mode) may resolve missingness for tabular static data, but such strategies may not be suitable for a variety of data types and architectures, ranging from static to temporal, tabular to imaging and sequencing data. As a result, advanced approaches are necessary to ensure model quality and robustness.

As described by Little and Rubin[3, 4], missing data can be categorized into three types: missing completely at random (MCAR), missing at random (MAR), and missing not at random (MNAR). Prior to the widespread adoption of deep learning (DL), traditional statistical and machine learning approaches, such as interpolation methods, k-nearest neighbor (k-NN)[5], multiple imputation by chained equations (MICE)[6], random forest (RF)-based models like MissForest[7], have been utilized to impute missing values. However, these widely used methods may be limited to certain categories of missing data; for instance, MICE generally requires the assumption that the data missing type is MAR[6]. When applied to complex healthcare data, these non-DL-based imputation strategies may also suffer from low accuracy[8, 9], especially when the mechanism of missingness is uncertain.

In recent years, DL-based methods have been used increasingly for solving missing value problems and demonstrated their ability in enhancing imputation accuracy[10, 11]. Furthermore, DL-based models can be tailored to handle complex missing patterns and data structures, such as time-series data with unique sequential structures and image data with spatial patterns[12, 13]. With superior performance and designation flexibility, DL-based imputation models have become popular in a variety of applications, from in-patient mortality prediction to early detection of Alzheimer's Disease (AD).

Although there are several existing reviews on missing value imputation, most of them either focus on non-DL-based methods[14, 15], or treat the neural network as a single type of method[16, 17]. Due to the lack of specificity, these articles cannot adequately assist prospective researchers who contemplate the application of DL-based models to their own data. To our knowledge, there has not been a systematic review of DL-based missing value imputation methods for diverse types of healthcare data. To address this gap, we present an evidence map analysis[18] that examines model-use by data type, and we aim to provide guidance for managing missing values with DL-based methodologies for researchers engaged in different clinical disciplines.

2. Materials and Methods

2.1 Search Strategies

In this study, we undertook a systematic literature search to identify relevant research articles. We searched five databases (MEDLINE, Web of Science, Embase, CINAHL, and Scopus) using a combination of search phrases "missing value", "imputation", "machine learning", "deep learning", and "healthcare". Detailed search strategies are provided in eTable 1.

2.2 Exclusion Criteria

We conducted the study according to the Preferred Reporting Items for Systematic Reviews (PRISMA) guidelines[19]. Papers with any of the following reasons were excluded: the study was not in the medical or clinical domain, the imputation model used was not DL-based or not specified, the study was not published as a research article (e.g., conference poster, conference abstract, or book chapter), or the article was not written in English.

2.3 Selection Procedure and Data Extraction

Two reviewers (ML and SL) independently screened the titles and abstracts based on the eligibility criteria between 6 August and 11 September 2021. Discrepancies were resolved through discussions with a third reviewer (HY). The full-text screening and information extraction took place from 12 September to 22 October 2021, during which ML and SL separately accessed the documents and, in the event of disagreement, consulted with HY. Four aspects of information were extracted from the included articles: data types, model backbones (i.e., main architectures), imputation strategies, and comparison with non-DL-based methods.

2.4 Data Analysis

To generate an evidence map[18] that illustrates the application of DL-based imputation models across various data types, we classified the types of data involved in imputation into six categories: tabular static data, tabular temporal data, genetic and genomic data, image data, signal data, and multi-modal data. Tabular static data and tabular temporal data both contain observations as rows and features as columns, but only tabular temporal data include the time factor. Genetic and genomic data encompasses both DNA data for organisms and personal genetic information. Image data and signal data refer to the information generated by specific medical devices, such as magnetic resonance imaging (MRI) and electrocardiogram (ECG). Multi-modal data refers to the use of multiple data types in a single imputation task.

Next, we categorized the articles according to the "backbones" of the imputation model: 1) multi-layer perceptron (MLP); 2) recurrent neural network (RNN), including vanilla RNN, long short-term memory (LSTM), and gated recurrent unit (GRU); 3) the framework of autoencoder (AE) which includes vanilla autoencoder, denoising autoencoder (DAE), and variational autoencoder (VAE); 4) the framework of generative adversarial network (GAN); 5) a hybrid of the four backbones mentioned above; and 6) other less common models, such as self-organizing map (SOM). The precise definitions of these models and their accompanying general imputation mechanisms are described in Table 1. To assess the imputation strategy, we divided it into two categories: separated and integrated. "Separated" means that the imputation process is separated from downstream tasks such as disease classification and risk prediction, whereas "integrated", also known as "end-to-end", means that the imputation process is addressed concurrently with downstream tasks.

The evidence map was presented based on the cross-tabulation of model backbones and data types. In addition, a bar plot displaying the distribution of imputation strategies was provided. Python version 3.8.3 (Python Software Foundation, Delaware, USA) and R version 4.0.2 (The R Foundation for Statistical Computing) were used for data analysis.

3. Results

Our search of five databases yielded 1058 papers, of which 64 were included for analysis. The details of the selection procedure are illustrated in Figure 1. Figure 2 depicts the evidence map between the "backbones" (i.e., main architectures) of DL-based missing value imputation models and the types of healthcare data. Among the 64 studies, 17 presented missing value imputation models for tabular static data, 24 for tabular temporal data, seven for genetic and genomic data, five for image data, five for signal data, and seven for multi-modal data. These numbers are not mutually exclusive as a single study may involve more than one data type. MLP model for tabular static data and vanilla autoencoder for tabular temporal data are the most common (9.4%, 6/64) among all cells in the cross-tabulation.

As indicated in Figure 3, most studies (70%, 45/64) adopted the "separated" strategy. The "integrated" strategy was popular among tabular temporal data and multi-modal data, but less used for tabular static data and genetic and genomic data, and rarely applied to image or signal data. Moreover, 35 of the 64 selected papers investigated the type of data missingness, 18 of which investigated MNAR missingness, while the rest developed imputation methods directly without examining this specification. Next, we present DL-based imputation techniques based on various types of health data.

3.1 Tabular Static Data

A total of 17 studies[20-36] used tabular static data in this review. Most (76%, 13/17) of them provided non-DL-based imputation methods as baselines for comparison with DL-based methods in terms of imputation accuracy (if complete data is available) and prediction performance, with simple imputation (mean/median/mode imputation, 41%, 7/17), RF-based methods (35%, 6/17), MICE (36%, 6/17), and k-NN (29%, 5/17) being the most common options. Most of these studies demonstrated the superiority of DL-based approaches.

The structures of the DL-based imputation methods for tabular static data are dominated by MLP-oriented models. Six studies used MLP models directly[23, 24, 30-32, 36], the other five created autoencoder models employing MLP modules as encoders and decoders[22, 25-27, 33], and three studies that utilized GANs also implemented simple MLP modules as generators and discriminators. In addition to MLP, this review also included several alternative options. Dynamic layered RNN[28] was applied to repeat the imputation process with an aim to enhance the imputation accuracy, an autoencoder with a structure of ladder network was developed to improve the imputation robustness in terms of different missing rates[34], and SOM was used in a comparison study which indicated the superiority of SOM in terms of improvement downstream task (classification) performance[29].

The handling of mixed variable types, such as numerical, categorical, and ordinal data, is a challenge for tabular data (both static and temporal), which is due to the assumptions about data distribution and limitations of some non-DL-based methods[26]. However, DL-based imputation methods are able to accommodate mixed variable types with proper encoding and activation functions[25].

3.2 Tabular Temporal Data

Among the 64 included studies, 24 [8, 9, 37-58] addressed missing value imputation for tabular temporal data. Missingness for this type of data is typically caused by factors that are less controllable, such as the dropping out of patients or the use of different assessment patterns for different patient subgroups[37]. Therefore, informative missingness with time dynamics and high missing rates pose challenges to imputation.

RNN-based (46%, 11/24) and AE-based (38%, 9/24) methods were commonly employed to impute tabular temporal data, where the "separated" imputation strategy was more often used in conjunction with the latter (27%, 3/11 for RNN-based versus 78%, 7/9 for AE-based). The development of the "separated" imputation strategy relied on the ground truth, that is, artificially composing missingness by dropping values at random and comparing the imputed values to the dropped ground truth. Among 64 studies, seven studies composed such missingness with various rates ranging from 5% to 90% in order to assess the robustness of imputation against different missing rates. Six (67%, 6/9) of the studies that applied AE-based methods artificially composed missing values and adopted the "separated" imputation strategy.

The nature of modeling sequence data makes RNN-based methods capable of capturing time dynamics-related missing patterns. Among these 11 studies that applied RNN-based methods, three studies[8, 49, 57] designed the models with bi (multi)-directions to utilize both past and future information for imputation, while another three[38, 40, 43] applied GRU-D – a model equipped with a specific parameter to characterize the decay of effects as time passes – or

developed variants of GRU-D. In addition, two studies developed hybrid RNN-based and AE-based methods in which the AE component was added following the initial imputation by the RNN component[39, 46]. Another two studies utilized the framework of GAN, where adversarial learning (alone[9] or in combination with RNN[45]) can contribute to avoiding error propagation from imputation to downstream tasks.

3.3 Genetic and Genomic Data

Seven studies[59-65] dealt with genetic and genomic data, including sequencing data[60, 62, 63, 65] and microarray data[59, 61, 64]. Sequencing data may comprise around 50% zero count observations[62, 65], some of which are “false” zeros, i.e., missing values due to insufficient sequencing input[60]. Furthermore, the high dimension property presents a barrier to the imputation of sequencing data. In microarray data, 90% of genes may have one or more missing values, as a consequence of experimental factors in the process, such as hybridization failures, artifacts on the microarray, insufficient resolution, image noise, and corruption, which are difficult to resolve[66].

Four studies[61-63, 65] used AE-based models (vanilla AE and VAE), and the remaining studies[59, 60, 64] employed MLP, where biological background knowledge may be incorporated into the imputation process for both types of models. With AE-based models, Badsha et al.[62] used transfer learning to extract prior knowledge about gene-gene relationships and learn the dependence structure within the reference panel. According to Kinalis et al.[63], the architecture of AE and VAE allows for more interpretable imputation procedures. With appropriate training, the autoencoder can be interpreted as a combination of biologically meaningful modules. Apart from AE’s natural ability to reduce the dimension in the latent space, MLP could also contribute to the development of the dimension deduction approach for imputation[59, 60]. As an example, Chen et al.[59] deciphered two low-dimensional hidden representations from the original high-dimensional data to explain a molecular relationship and a sample level connection, respectively.

3.4 Image Data

Among the five papers[12, 67-70] that focused on image data, five targeted neuroimages, including magnetic resonance imaging (MRI) and positron emission tomography (PET) that are widely used for computer-aided diagnosis of AD and mild cognitive impairment (MCI)[12, 68-70]. Another study[67] dealt with cardiac magnetic resonance (CMR) images, which are often considered the gold standard for many cardiovascular medicine analyses. Insufficient image quality and acquisition or storage errors are common causes of pixel-level missingness in image data[12, 67]. In practice, MRI patients may reject PET scans due to their high cost and radioactive exposure[70], resulting in the absence of the entire image.

All five studies designed GAN frameworks for missing value imputation. Only one study[67] established non-DL-based imputation baselines for comparison (mean imputation and an interpolation approach), while the other four did not. Conditional GAN that utilizes label information could be used to impute missing pixels in an image or to synthesize a new one when the whole image is missing[67]. Imputing PET data based on the corresponding MRI data is natural for diagnosing brain diseases such as AD since the usage of multi-source images can provide rich information[12, 71]. In this case, the entire image rather than a few pixel values must be imputed. To impute the entire image, a task-induced GAN can be developed to improve the imputation quality with two tasks designed for the discriminator: whether the image is true (the imputed MRI must be true enough to fool the discriminator) and whether it indicates the disease (the imputed MRI should be consistent with PET)[12]. Cycle-GAN and Colla-GAN have been adopted for this task, which required two discriminators to detect synthetic MRI and PET, respectively[68, 69], similar to the approach proposed by Pan et al. that utilized the connection between MRI and PET and designed two generators and two discriminators in their GAN framework[70].

3.5 Signal Data

In this review, various signal data were discussed in five studies, including actigraphy device data[72], smartphone applications data[73], wearable sensors data[74], medical waveform (e.g.

ECG) data[75], time-series measured at each voxel of fMRI[76]. Compared with tabular temporal data, signal data usually has a high sample rate and is susceptible to noise. For data collected in the movement, a common issue is the large amount of missing data with blocks of features concurrently lost as a result of intermittent disconnections, body movement, and firmware malfunction[74], where the missing interval may need to be characterized as a period of continuously repeated zero values[72]. For medical waveform data, even a low missing rate in real-time operational systems, such as removing the lead for a few seconds, can significantly harm the prediction performance over the next hour[75]. Due to the fMRI technique's great susceptibility to noise artifacts, as with image modality, missing values become common for signal modality as well, posing challenges to the imputation process, where the spatiotemporal nature can be of assistance[76].

The non-DL-based imputation method of simple imputation (mean/median imputation, 80%, 4/5) is a common baseline in these five studies. Three studies used AE-based models with customization, such as denoising autoencoder where missingness was treated as a type of noise[72], adversarial autoencoder where the encoder contributed to feature representation and followed by the discriminator of GAN[73], and 1-D convolutional and corresponding deconvolutional modules used by Miller et al. in their autoencoder framework[75]. The other two studies[74, 76] applied RNN models with the time factor into consideration. For signal imputation in fMRI data, the missing values were first filled based on spatial information, and a GRU layer was employed for regularization in the time domain[76].

3.6 Multi-modal Data

Seven studies[10, 77-82] investigated imputation methods for multi-modal data, where information fusion procedures were essential in connecting modality-specific models. Four employed autoencoder models and designed linkages between encoders and decoders to concatenate different modalities[10, 77, 78, 82]. For example, Kim et al.[78] developed a stacked DAE model with a merged hidden layer working as the linkage. Additionally, Kim et al.[79] created a collaborative layer to connect MLP models for different modalities. The task-specific

layers designed by Thung et al.[80] could enable iterative communication between modality-specific layers, thereby facilitating the exchange of cross-modal information. Vivar et al.[81] established an end-to-end framework, in which the self-attention process aggregated recurrent graph convolutional models. Using this architecture, missing value imputation was transformed into the completion of a geometric matrix. Combining the graph convolutional model with LSTM has shown to be effective in solving this geometric problem.

4. Discussion

In this systematic review, we contribute to exhaustively summarizing knowledge regarding the efficacy of deep learning models by investigating their applications to missing value imputation for healthcare data. Having identified popular deep learning models with respect to data type, we found that deep learning models could be customized to take both the data type and the corresponding missing patterns into consideration, thereby improving the quality of missing value imputation. Moreover, the “integrated” imputation strategy could enhance the performance of both imputation and downstream analysis, and the usage varied across data types, highlighting the benefit of imputation based on the data type’s characteristics.

Designing and implementing data-type-oriented imputation models could benefit both the imputation process and downstream tasks. The MLP models and autoencoders that use MLP as encoders and decoders are well-suited to capturing the feature relevance of tabular static data in a relatively straightforward manner. When dealing with tabular temporal data, informative missingness and high missing rates are common[38], which makes it difficult to characterize time dynamics. RNN-based models such as GRU-D[38] and its variants are helpful for capturing time patterns and imputing missing values for tabular temporal data. For genetic and genomic data, MLP or AE with MLP acting as encoder and decoder is appropriate for imputation with effectively used biological knowledge.

GAN-based models are commonly used for image imputation. Additional information, such as labeling (e.g., used in Co-GAN[83]) and relevant images (as used in Cycle-GAN[84] or CollaGAN[69]), can enhance the performance of imputation. Modules from CNN to Transformer[85] to Swin Transformer[86] could be added to the framework to tackle spatial information in image data. The imputation procedure for general signal data may resemble that for tabular temporal data. In multi-model data, the fusion of mode-specific models is indispensable, and most of the current operations are focused at the layer level, for example, stacking and the self-attention mechanism.

There is currently a vacancy in the imputation approach for medical text data. This may be explained by the fact that techniques in natural language processing (e.g., BERT[87]) inherently learn representation through masking, that is, considering some language tokens as missing on purpose, so actual missingness will not be an issue. Healthcare text data can be analyzed using customized biomedical research models, such as BioBERT[88] and MedBERT[89].

Due to the widespread use of diverse and complex data types, the causes and forms of missing data have become more complicated. Instead of imputing data of different types separately, hybrid models are designed to better capture missing patterns and benefit from the interaction among imputation processes through an integral approach[78, 79]. For example, the integration of LSTM and the graphical convolutional model developed by Vivar et al.[81] advanced the performance of imputation and classification in computer-aided diagnosis. Developing hybrid models capable of handling hybrid data will open up possibilities for further research on imputation in complex scenarios. Rather than referring to mixed data types, some researchers use the term "multi-modality" to describe datasets from different sources but of the same type (e.g., MRI and PET images[12, 68]). For this image type of "multi-modality", GAN-based frameworks that involve image translation[69] and Transformer-based architectures that are intrinsically integrated with attention modules may also be able to solve the imputation problem[85].

Block-building logic enables DL-based models to adopt an “integrated” strategy, i.e., co-training imputation and downstream tasks, which is preferred for several reasons. First, the interaction between these two tasks can be mutually beneficial, reducing the bias in imputation, and providing prior information for the downstream modelling[39, 45, 81]. Second, the “integrated” strategy could be more practical because it circumvents the difficulties in defining imputation accuracy when the missing rate is high, indicating limited ground truth for imputation quality checking[90], or when multiple data types are involved and the overall imputation accuracy is difficult to define. As shown in Figure 3, this is consistent with the relative prevalence of the “integrated” strategy when working with tabular temporal and multi-modal data.

Third, in contrast to the “separated” strategy, the “integrated” strategy does not emphasize the selection of the optimal combination of imputation and downstream models[81]. The step-wise selection employed in the “separated” strategy, i.e., determining the best imputation model first and then deciding on the downstream models, may not be effective given the belief that the accuracy of imputation does not directly affect the accuracy of downstream tasks[24, 81]. Although it is theoretically possible to test all potential combinations of imputation and downstream models and select the optimal combination, doing so would be extremely time-consuming. The “integrated” strategy resolves such practical difficulties by imputing missing data while developing the downstream models. Last, the “integrated” strategy requires little additional effort when both the imputation and downstream models are DL-based, since the parameters can easily be linked up and updated simultaneously. However, the additional model complexity resulting from the “integrated” strategy can restrict its wide application, which explains its limited use in current studies (see Figure 3).

When both non-DL-based and DL-based imputation models are available, the former may be preferred for its ease of implementation and the fact that non-DL-based methods (such as MICE, Xgboost, LightGBM, etc.) can produce good imputation results when coupled with carefully engineered tabular static and temporal data[14]. The ease of implementation is, however, dependent upon restrictive statistical assumptions about data, which are usually difficult to

examine and identify in real-world scenarios[91]. The effort required for feature engineering could be tremendous[27], thereby distracting researchers from their primary research objectives. Moreover, other concerns, such as high dimension[55, 58] and low time efficiency[55, 78], pose obstacles for non-DL-based methods (e.g., k-NN, MICE, etc.). In the presence of healthcare data in complex formats, DL-based models may be a good solution, as statistical assumptions and feature engineering are rarely required and they do not suffer from the curse of dimensionality as many non-DL-based methods do. Additionally, since DL-based models can be directly pre-trained and applied to new samples[61], they may reduce computational costs during evaluation.

Clinical practitioners who lack expertise in deep learning may find it difficult to implement DL-based imputation models, so DL-based models that are specified by data type and packaged well are of practical value. Moreover, researchers should also pay attention to fairness in the imputation process, which does not receive adequate attention at this time, as most studies in this review refer to statistical bias, but do not address social bias, or discrimination against certain groups or individuals[1, 13]. Data imputation influenced by such bias may degrade the subsequent analysis and lead to unfair decision-making, thereby causing inequality and burdening the healthcare system.

This study has several limitations. First, the scope of our review was limited to clinical and translational research; however, alternative DL-based imputation techniques may exist in other research fields and merit examination in a future study. Second, we concentrated primarily on data types and their corresponding imputation strategies, while a more in-depth analysis based on specific research topics may yield additional insights. A discussion of healthcare applications that goes deeper by examining each paper's particular research topics may contribute to the development of effective imputation models.

5. Conclusions

Our study fills a gap in the existing literature where DL-based missing value imputation methods have not been systematically reviewed and evaluated. We thoroughly reviewed the DL-based imputation models across a variety of healthcare data types and summarized the missingness mechanisms and appropriate models for each data type. In addition, we identified the advantages of DL-based imputation models and the corresponding integrated strategy. Future research may focus on the portability and fairness of DL-based imputation models with applications to healthcare data.

Reference

- [1] R. B. Parikh, S. Teeple, and A. S. Navathe, "Addressing Bias in Artificial Intelligence in Health Care," *JAMA*, vol. 322, no. 24, pp. 2377-2378, 2019, doi: 10.1001/jama.2019.18058.
- [2] H. Kang, "The prevention and handling of the missing data," (in eng), *Korean journal of anesthesiology*, vol. 64, no. 5, pp. 402-406, 2013, doi: 10.4097/kjae.2013.64.5.402.
- [3] R. Little, R. J. A. Little, and D. B. Rubin, *Statistical analysis with missing data*. Wiley, 2019.
- [4] D. B. Rubin, "Inference and missing data," *Biometrika*, vol. 63, no. 3, pp. 581-592, 1976, doi: 10.1093/biomet/63.3.581.
- [5] G. E. A. P. A. Batista and M. C. Monard, "An analysis of four missing data treatment methods for supervised learning," *Applied Artificial Intelligence*, vol. 17, no. 5-6, pp. 519-533, 2003/05/01 2003, doi: 10.1080/713827181.
- [6] S. van Buuren and K. Groothuis-Oudshoorn, "mice: Multivariate Imputation by Chained Equations in R," *Journal of Statistical Software*, vol. 45, no. 3, pp. 1 - 67, 12/12 2011, doi: 10.18637/jss.v045.i03.
- [7] D. J. Stekhoven and P. Bühlmann, "MissForest—non-parametric missing value imputation for mixed-type data," *Bioinformatics*, vol. 28, no. 1, pp. 112-118, 2011, doi: 10.1093/bioinformatics/btr597.
- [8] J. Yoon, W. R. Zame, M. van der Schaar, J. Yoon, W. R. Zame, and M. van der Schaar, "Estimating Missing Data in Temporal Data Streams Using Multi-Directional Recurrent

- Neural Networks," *IEEE TRANSACTIONS ON BIOMEDICAL ENGINEERING*, vol. 66, no. 5, pp. 1477-1490, 2019.
- [9] S. Park *et al.*, "Learning Sleep Quality from Daily Logs," *KDD'19: PROCEEDINGS OF THE 25TH ACM SIGKDD INTERNATIONAL CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING*, pp. 2421-2429, 2019.
- [10] M. Peralta, P. Jannin, C. Haegelen, and J. S. H. Baxter, "Data imputation and compression for Parkinson's disease clinical questionnaires," *Artificial Intelligence in Medicine*, vol. 114, p. 102051, 2021/04/01/ 2021, doi: <https://doi.org/10.1016/j.artmed.2021.102051>.
- [11] D. Xu *et al.*, "A deep learning-based, unsupervised method to impute missing values in electronic health records for improved patient management," (in English), *JOURNAL OF BIOMEDICAL INFORMATICS*, vol. 111, 2020-11 2020.
- [12] X. Gao, F. Shi, D. Shen, and M. Liu, "Task-induced Pyramid and Attention GAN for Multimodal Brain Image Imputation and Classification in Alzheimers disease," *IEEE journal of biomedical and health informatics*, 2021. [Online]. Available: <https://pubmed.ncbi.nlm.nih.gov/34280112/>.
- [13] D. Xu *et al.*, "A Deep Learning-Based Unsupervised Method to Impute Missing Values in Patient Records for Improved Management of Cardiovascular Patients," (in English), *IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS*, vol. 25, no. 6, pp. 2260-2272, 2021-6 2021.
- [14] Y. Luo, "Evaluating the state of the art in missing data imputation for clinical data," *Briefings in Bioinformatics*, 2021, doi: 10.1093/bib/bbab489.
- [15] J. C. Jakobsen, C. Glud, J. Wetterslev, and P. Winkel, "When and how should multiple imputation be used for handling missing data in randomised clinical trials – a practical guide with flowcharts," *BMC Medical Research Methodology*, vol. 17, no. 1, p. 162, 2017/12/06 2017, doi: 10.1186/s12874-017-0442-1.
- [16] T. Thomas and E. Rajabi, "A systematic review of machine learning-based missing value imputation techniques," *Data Technologies and Applications*, vol. 55, no. 4, pp. 558-585, 2021, doi: 10.1108/DTA-12-2020-0298.

- [17] S. Jäger, A. Allhorn, and F. Bießmann, "A Benchmark for Data Imputation Methods," (in English), *Frontiers in Big Data*, Original Research vol. 4, 2021-July-08 2021, doi: 10.3389/fdata.2021.693674.
- [18] N. Liu *et al.*, "Coronavirus disease 2019 (COVID-19): an evidence map of medical literature," (in eng), *BMC Med Res Methodol*, vol. 20, no. 1, p. 177, Jul 2 2020, doi: 10.1186/s12874-020-01059-y.
- [19] A. Liberati *et al.*, "The PRISMA statement for reporting systematic reviews and meta-analyses of studies that evaluate healthcare interventions: explanation and elaboration," *BMJ*, vol. 339, p. b2700, 2009, doi: 10.1136/bmj.b2700.
- [20] W. Dong *et al.*, "Generative adversarial networks for imputing missing data for big data clinical research," *BMC Med Res Methodol*, vol. 21, no. 1, p. 78, Apr 20 2021, doi: 10.1186/s12874-021-01272-3.
- [21] E. Hallaji *et al.*, "Adversarial Learning on Incomplete and Imbalanced Medical Data for Robust Survival Prediction of Liver Transplant Patients," (in English), *IEEE ACCESS*, vol. 9, pp. 73641-73650, 2021 2021.
- [22] E. Macias *et al.*, "Novel Imputation Method Using Average Code from Autoencoders in Clinical Data," *28TH EUROPEAN SIGNAL PROCESSING CONFERENCE (EUSIPCO 2020)*, pp. 1576-1579, 2021.
- [23] C. Y. Cheng *et al.*, "A Deep Learning Approach for Missing Data Imputation of Rating Scales Assessing Attention-Deficit Hyperactivity Disorder," *FRONTIERS IN PSYCHIATRY*, vol. 11, 2020.
- [24] D. Vrbaski *et al.*, "Missing Data Imputation in Cardiometabolic Risk Assessment: A Solution Based on Artificial Neural Networks," *COMPUTER SCIENCE AND INFORMATION SYSTEMS*, vol. 17, no. 2, pp. 379-401, 2020.
- [25] N. Abiri *et al.*, "Establishing strong imputation performance of a denoising autoencoder in a wide range of missing data problems," (in English), *NEUROCOMPUTING*, vol. 365, pp. 137-146, 2019-11-6 2019.
- [26] K. Miok *et al.*, "Multiple Imputation for Biomedical Data using Monte Carlo Dropout Autoencoders," *2019 E-HEALTH AND BIOENGINEERING CONFERENCE (EHB)*, 2019.

- [27] S. Phung, A. Kumar, J. Kim, S. Phung, A. Kumar, and J. G. P. I. E. E. Kim, "A deep learning technique for imputing missing healthcare data," *2019 41ST ANNUAL INTERNATIONAL CONFERENCE OF THE IEEE ENGINEERING IN MEDICINE AND BIOLOGY SOCIETY (EMBC)*, pp. 6513-6516, 2019.
- [28] H. Turabieh, A. Abu Salem, N. Abu-El-Rub, H. Turabieh, A. Abu Salem, and N. Abu-El-Rub, "Dynamic L-RNN recovery of missing data in IoMT applications," *FUTURE GENERATION COMPUTER SYSTEMS-THE INTERNATIONAL JOURNAL OF ESCIENCE*, vol. 89, pp. 575-583, 2018.
- [29] E. M. Hernandez-Pereira, D. Alvarez-Estevez, V. Moret-Bonillo, E. M. Hernandez-Pereira, D. Alvarez-Estevez, and V. Moret-Bonillo, "Automatic classification of respiratory patterns involving missing data imputation techniques," *BIOSYSTEMS ENGINEERING*, vol. 138, pp. 65-76, 2015.
- [30] M. W. Huang, W. C. Lin, and C. F. Tsai, "Outlier Removal in Model-Based Missing Value Imputation for Medical Datasets," *Journal of healthcare engineering*, vol. 2018, pp. 1817479-1817479, 2018. [Online]. Available: <https://pubmed.ncbi.nlm.nih.gov/29599943/>.
- [31] W. Seffens, C. Evans, and H. Taylor, "Machine Learning Data Imputation and Classification in a Multicohort Hypertension Clinical Study," *Bioinformatics and biology insights*, vol. 9, pp. 43-54, 2015. [Online]. Available: <https://pubmed.ncbi.nlm.nih.gov/27199552/>.
- [32] C. M. Ennett, M. Frize, and C. Walker, "Imputation of missing values by integrating neural networks and case-based reasoning," *Annual International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Annual International Conference*, vol. 2008, pp. 4337-41, 2008. [Online]. Available: <https://pubmed.ncbi.nlm.nih.gov/19163673/>.
- [33] Z. Chen, B. Cao, A. Edwards, H. Deng, and K. Zhang, "A deep imputation and inference framework for estimating personalized and race-specific causal effects of genomic alterations on PSA," *J Bioinform Comput Biol*, pp. 2150016-2150016, 2021. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L635542803&from=export> U2 - L635542803.

- [34] E. Hallaji, R. Razavi-Far, and M. Saif, "DLIN: Deep Ladder Imputation Network," (in English), *IEEE Trans Cybern*, 2021 2021. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L634516500&from=export> U2 - L634516500.
- [35] M. Kachuee, K. Karkkainen, O. Goldstein, S. Darabi, and M. Sarrafzadeh, "Generative Imputation and Stochastic Prediction," (in English), *IEEE Trans Pattern Anal Mach Intell*, 2020 2020. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L632815301&from=export> U2 - L632815301.
- [36] J. Bektaş, T. Ibrikçi, and İ. T. Özcan, "The impact of imputation procedures with machine learning methods on the performance of classifiers: An application to coronary artery disease data including missing values," *Biomed. Res.*, vol. 29, no. 13, pp. 2780-2785, 2018. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L623713839&from=export> U2 - L623713839.
- [37] M. M. Ghazi *et al.*, "Training recurrent neural networks robust to incomplete data: Application to Alzheimer's disease progression modeling," *MEDICAL IMAGE ANALYSIS*, vol. 53, pp. 39-46, 2019.
- [38] Z. P. Che *et al.*, "Recurrent Neural Networks for Multivariate Time Series with Missing Values," *SCIENTIFIC REPORTS*, vol. 8, 2018.
- [39] J. de Jong *et al.*, "Deep learning for clustering of multivariate clinical patient trajectories with missing values," *GIGASCIENCE*, vol. 8, no. 11, 2019.
- [40] M. Habiba and B. A. Pearlmutter, "Neural ODEs for Informative Missingness in Multivariate Time Series," *2020 31st Irish Signals and Systems Conference (ISSC)*, pp. 1-6, 2020.
- [41] W. Jung *et al.*, "Deep recurrent model for individualized prediction of Alzheimer's disease progression," *NEUROIMAGE*, vol. 237, 2021.
- [42] T. Tsiligkaridis and J. Sloboda, "A Multi-task LSTM Framework for Improved Early Sepsis Prediction," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 12299, pp. 49-58, 2020. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0->

- 85092219124&doi=10.1007%2f978-3-030-59137-3_5&partnerID=40&md5=bfcf4279d270db1d310221765f6c1742.
- [43] Y. Zhao, M. Berretta, T. Wang, and T. Chitnis, "GRU-DF: A Temporal Model with Dynamic Imputation for Missing Target Values in Longitudinal Patient Data," *2020 IEEE International Conference on Healthcare Informatics, ICHI 2020*, 2020. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85103205663&doi=10.1109%2fICHI48887.2020.9374359&partnerID=40&md5=f52af026aa4b7815539b181a21720a53>.
- [44] W. Jung, A. W. Mulyadi, H. I. Suk, W. Jung, A. W. Mulyadi, and H.-I. Suk, "Unified Modeling of Imputation, Forecasting, and Prediction for AD Progression," *MEDICAL IMAGE COMPUTING AND COMPUTER ASSISTED INTERVENTION - MICCAI 2019, PT IV*, vol. 11767, pp. 168-176, 2019.
- [45] Q. Ma, S. Li, and G. Cottrell, "Adversarial Joint-Learning Recurrent Neural Network for Incomplete Time Series Classification," (in English), *IEEE Trans Pattern Anal Mach Intell*, 2020 2020. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L633082401&from=export U2 - L633082401>.
- [46] A. W. Mulyadi, E. Jun, and H. Suk, "Uncertainty-Aware Variational-Recurrent Imputation Network for Clinical Time Series," *IEEE Transactions on Cybernetics*, 2021. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85102286823&doi=10.1109%2fTCYB.2021.3053599&partnerID=40&md5=7081d5e4fa8c5b77c692654e7e67a77>.
- [47] B. K. Beaulieu-Jones and J. H. Moore, "MISSING DATA IMPUTATION IN THE ELECTRONIC HEALTH RECORD USING DEEPLY LEARNED AUTOENCODERS," *Pac Symp Biocomput*, vol. 22, pp. 207-218, 2016. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L617085279&from=export U2 - L617085279>.
- [48] F. M. Bianchi *et al.*, "Learning representations of multivariate time series with missing data," *PATTERN RECOGNITION*, vol. 96, 2019.

- [49] J. Codella, H. Sarker, P. Chakraborty, M. Ghalwash, Z. Yao, and D. Sow, "EXITs: An Ensemble Approach for Imputing Missing EHR Data," *2019 IEEE International Conference on Healthcare Informatics, ICHI 2019*, 2019. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85075913744&doi=10.1109%2fICHI.2019.8904779&partnerID=40&md5=4c1cf24bce38c234d5e1aafd197cdbe8>.
- [50] Y. Ding, W. N. Street, L. Tong, and S. Wang, "An ensemble method for data imputation," *2019 IEEE International Conference on Healthcare Informatics, ICHI 2019*, 2019. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85075923066&doi=10.1109%2fICHI.2019.8904629&partnerID=40&md5=f5b98e17391f1d306fed8c1e7651c16b>.
- [51] V. Fortuin *et al.*, "GP-VAE: Deep Probabilistic Multivariate Time Series Imputation," *INTERNATIONAL CONFERENCE ON ARTIFICIAL INTELLIGENCE AND STATISTICS, VOL 108*, vol. 108, pp. 1651-1660, 2020.
- [52] E. Jun, A. W. Mulyadi, J. Choi, and H. Suk, "Uncertainty-Gated Stochastic Sequential Model for EHR Mortality Prediction," *IEEE Transactions on Neural Networks and Learning Systems*, 2020. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85091316861&doi=10.1109%2fTNNLS.2020.3016670&partnerID=40&md5=c0273ef5abb40659d056f281a06b7ce5>.
- [53] E. Jun, A. W. Mulyadi, and H. I. Suk, "Stochastic Imputation and Uncertainty-Aware Attention to EHR for Mortality Prediction," in *2019 International Joint Conference on Neural Networks (IJCNN)*, 14-19 July 2019 2019, pp. 1-7, doi: 10.1109/IJCNN.2019.8852132.
- [54] S. W. Lin *et al.*, "Filling Missing Values on Wearable-Sensory Time Series Data," *PROCEEDINGS OF THE 2020 SIAM INTERNATIONAL CONFERENCE ON DATA MINING (SDM)*, pp. 46-54, 2020.
- [55] D. Xu *et al.*, "A deep learning-based, unsupervised method to impute missing values in electronic health records for improved patient management," *JOURNAL OF BIOMEDICAL INFORMATICS*, vol. 111, 2020.

- [56] D. Xu *et al.*, "A Deep Learning-Based Unsupervised Method to Impute Missing Values in Patient Records for Improved Management of Cardiovascular Patients," *IEEE JOURNAL OF BIOMEDICAL AND HEALTH INFORMATICS*, vol. 25, no. 6, pp. 2260-2272, 2021.
- [57] C. Yin, R. Liu, D. Zhang, and P. Zhang, "Identifying Sepsis Subphenotypes via Time-Aware Multi-Modal Auto-Encoder," *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 862-872, 2020. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85090419687&doi=10.1145%2f3394486.3403129&partnerID=40&md5=6e3023afa1a3fc08c3f178ad0185642a>.
- [58] H. M. Tao, Q. Q. Deng, S. Z. Xiao, T. Huamin, D. Qiuqun, and X. Shanzhu, "Reconstruction of time series with missing value using 2D representation-based denoising autoencoder," *JOURNAL OF SYSTEMS ENGINEERING AND ELECTRONICS*, vol. 31, no. 6, pp. 1087-1096, 2020.
- [59] L. Chen, J. Xu, and S. C. Li, "DeepMF: deciphering the latent patterns in omics profiles with a deep learning method," *BMC Bioinformatics*, vol. 20, pp. 648-648, 2019. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L630464729&from=export U2 - L630464729>.
- [60] A. Mongia, D. Sengupta, and A. Majumdar, "deepMc: Deep Matrix Completion for Imputation of Single-Cell RNA-seq Data," (in English), *J. Comput. Biol.*, vol. 27, no. 7, pp. 1011-1019, 2020 2020. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L629708689&from=export U2 - L629708689>.
- [61] Y. L. Qiu, H. Zheng, and O. Gevaert, "Genomic data imputation with variational auto-encoders," *GigaScience*, vol. 9, no. 8, 2020. [Online]. Available: <https://pubmed.ncbi.nlm.nih.gov/32761097/>.
- [62] M. B. Badsha *et al.*, "Imputation of single-cell gene expression with an autoencoder neural network," *QUANTITATIVE BIOLOGY*, vol. 8, no. 1, pp. 78-94, 2020.
- [63] S. Kinalis, F. C. Nielsen, O. Winther, and F. O. Bagger, "Deconvolution of autoencoders to learn biological regulatory modules from single cell mRNA sequencing data," *BMC*

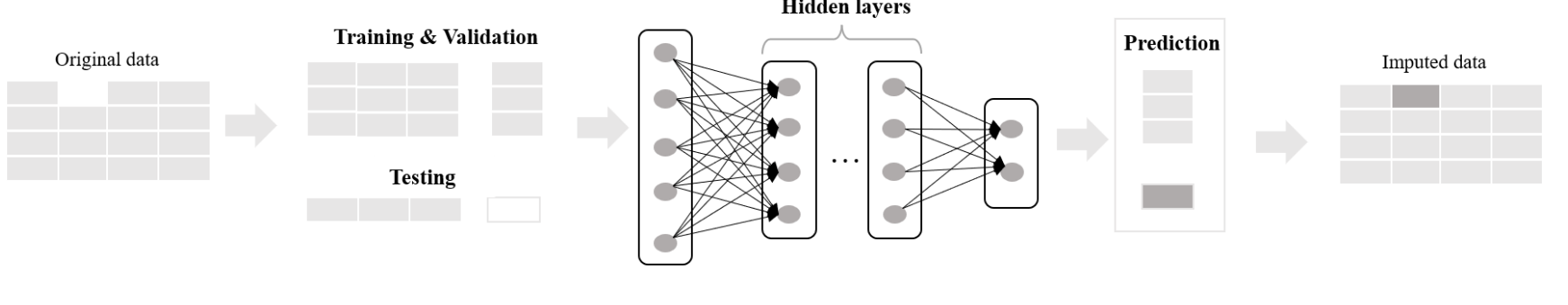
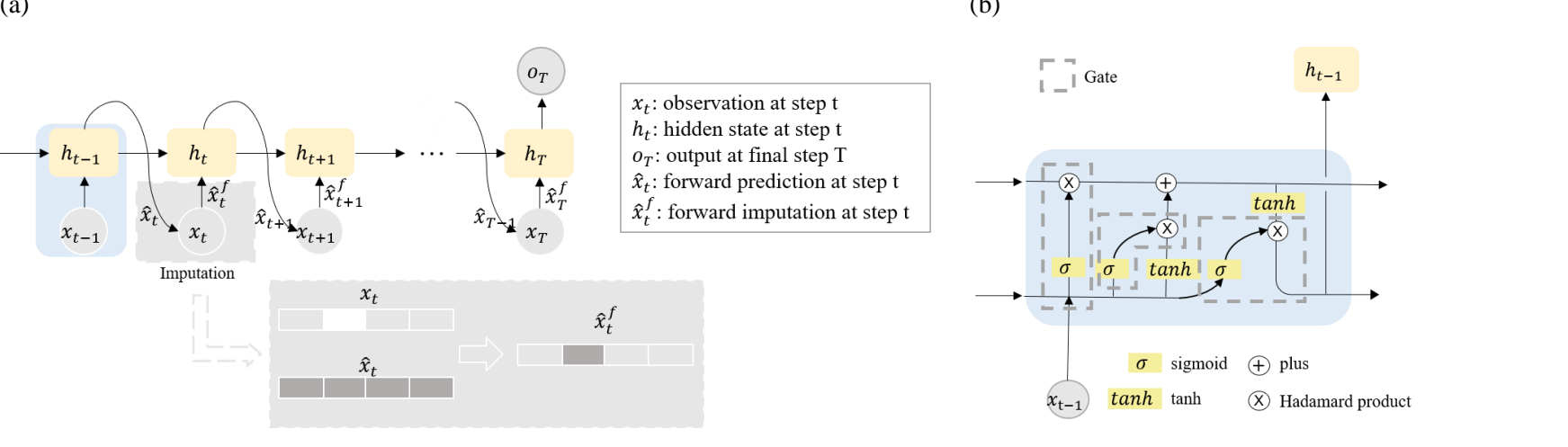
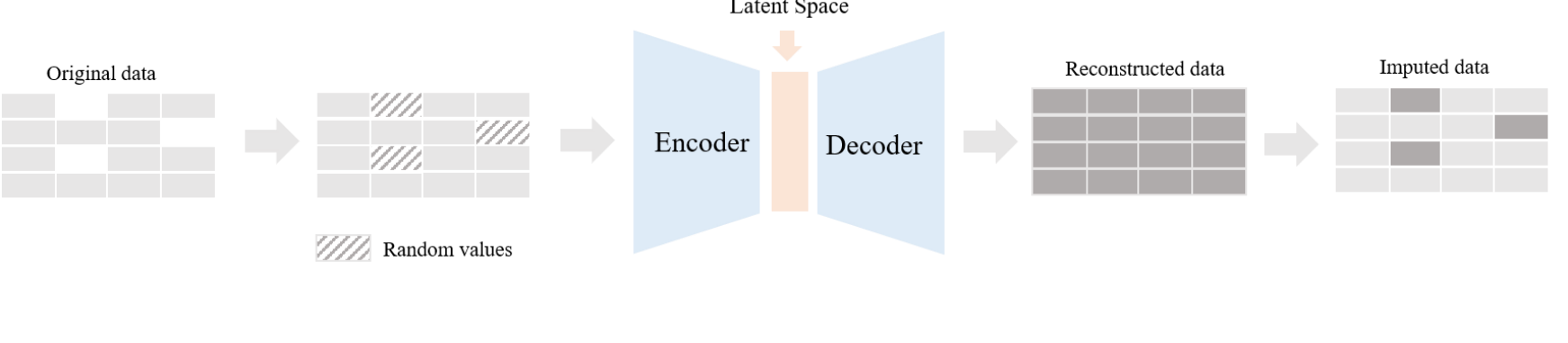
- bioinformatics*, vol. 20, no. 1, pp. 379-379, 2019. [Online]. Available: <https://pubmed.ncbi.nlm.nih.gov/31286861/>.
- [64] Y. V. Sun and S. L. R. Kardia, "Imputing missing genotypic data of single-nucleotide polymorphisms using neural networks," *Eur. J. Hum. Genet.*, vol. 16, no. 4, pp. 487-495, 2008. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L351420382&from=export> U2 - L351420382.
- [65] T. Tian, M. R. Min, Z. Wei, T. Tian, M. R. Min, and Z. Wei, "Model-based autoencoders for imputing discrete single-cell RNA-seq data," *METHODS*, vol. 192, 2021.
- [66] A. W.-C. Liew, N.-F. Law, and H. Yan, "Missing value imputation for gene expression data: computational techniques to recover missing data from available information," *Briefings in Bioinformatics*, <https://doi.org/10.1093/bib/bbq080> vol. 12, no. 5, pp. 498-513, 2010. [Online]. Available: <https://doi.org/10.1093/bib/bbq080>.
- [67] Y. Xia *et al.*, "Recovering from missing data in population imaging - Cardiac MR image imputation via conditional generative adversarial nets," *MEDICAL IMAGE ANALYSIS*, vol. 67, 2021.
- [68] Y. S. Pan *et al.*, "Synthesizing Missing PET from MRI with Cycle-consistent Generative Adversarial Networks for Alzheimer's Disease Diagnosis," *MEDICAL IMAGE COMPUTING AND COMPUTER ASSISTED INTERVENTION, PT III*, vol. 11072, pp. 455-463, 2018.
- [69] D. Lee, J. Kim, W. J. Moon, and J. C. Ye, "Collagan: Collaborative gan for missing image data imputation," *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2019, pp. 2482-2491, 2019. [Online]. Available: <https://www.scopus.com/inward/record.uri?eid=2-s2.0-85077542195&doi=10.1109%2fCVPR.2019.00259&partnerID=40&md5=57709225be62b798bb943e5ebf79ad04>.
- [70] Y. S. Pan *et al.*, "Spatially-Constrained Fisher Representation for Brain Disease Identification With Incomplete Multi-Modal Neuroimages," *IEEE TRANSACTIONS ON MEDICAL IMAGING*, vol. 39, no. 9, pp. 2965-2975, 2020.
- [71] B. D. James, S. E. Leurgans, L. E. Hebert, P. A. Scherr, K. Yaffe, and D. A. Bennett, "Contribution of Alzheimer disease to mortality in the United States," (in eng),

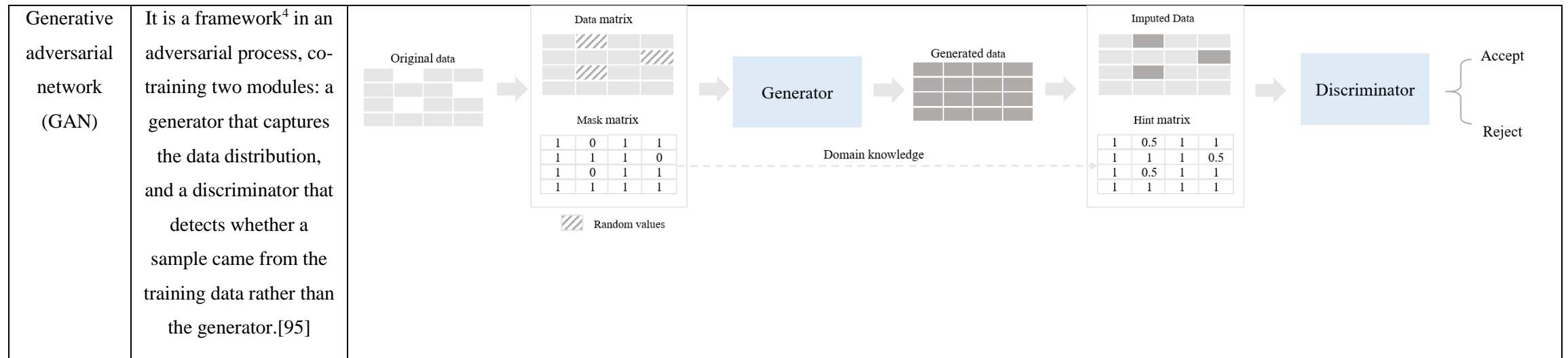
- Neurology*, vol. 82, no. 12, pp. 1045-50, Mar 25 2014, doi: 10.1212/wnl.0000000000000240.
- [72] J. H. Jang *et al.*, "Deep Learning Approach for Imputation of Missing Values in Actigraphy Data: Algorithm Development Study," *JMIR MHEALTH AND UHEALTH*, vol. 8, no. 7, 2020.
- [73] A. Saeed, T. Ozcelebi, and J. Lukkien, "Synthesizing and Reconstructing Missing Sensory Modalities in Behavioral Context Recognition," (in English), *Sensors (Basel)*, vol. 18, no. 9, 2018 2018. [Online]. Available: <https://www.embase.com/search/results?subaction=viewrecord&id=L624350955&from=export U2 - L624350955>.
- [74] T. Feng and S. Narayanan, "Imputing Missing Data In Large-Scale Multivariate Biomedical Wearable Recordings Using Bidirectional Recurrent Neural Networks With Temporal Activation Regularization," in *2019 41st Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC)*, 23-27 July 2019 2019, pp. 2529-2534, doi: 10.1109/EMBC.2019.8856966.
- [75] D. Miller *et al.*, "Physiological Waveform Imputation of Missing Data using Convolutional Autoencoders," *2018 IEEE 20TH INTERNATIONAL CONFERENCE ON E-HEALTH NETWORKING, APPLICATIONS AND SERVICES (HEALTHCOM)*, 2018.
- [76] D. Calhas and R. Henriques, "fMRI Multiple Missing Values Imputation Regularized by a Recurrent Denoiser," *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 12721, pp. 25-35, 2021. [Online]. Available: https://www.scopus.com/inward/record.uri?eid=2-s2.0-85111355867&doi=10.1007%2f978-3-030-77211-6_3&partnerID=40&md5=a66e6f24ae4d54d82ef27ab646850bde.
- [77] H. Akramifard, M. A. Balafar, S. N. Razavi, and A. R. Ramli, "Early Detection of Alzheimer's Disease Based on Clinical Trials, Three-Dimensional Imaging Data, and Personal Information Using Autoencoders," (in eng), *J Med Signals Sens*, vol. 11, no. 2, pp. 120-130, Apr-Jun 2021, doi: 10.4103/jmss.JMSS_11_20.
- [78] J. C. Kim, K. Chung, J.-C. Kim, and K. Chung, "Multi-Modal Stacked Denoising Autoencoder for Handling Missing Data in Healthcare Big Data," *IEEE ACCESS*, vol. 8, pp. 104933-104943, 2020.

- [79] J. C. Kim, K. Chung, J.-C. Kim, and K. Chung, "Hybrid Multi-Modal Deep Learning using Collaborative Concat Layer in Health Bigdata," *IEEE ACCESS*, vol. 8, pp. 192469-192480, 2020.
- [80] K. H. Thung, P. T. Yap, D. G. Shen, K.-H. Thung, P.-T. Yap, and D. Shen, "Multi-stage Diagnosis of Alzheimer's Disease with Incomplete Multimodal Data via Multi-task Deep Learning," *DEEP LEARNING IN MEDICAL IMAGE ANALYSIS AND MULTIMODAL LEARNING FOR CLINICAL DECISION SUPPORT*, vol. 10553, pp. 160-168, 2017.
- [81] G. Vivar *et al.*, "Simultaneous imputation and classification using Multigraph Geometric Matrix Completion (MGMC): Application to neurodegenerative disease classification," *ARTIFICIAL INTELLIGENCE IN MEDICINE*, vol. 117, 2021.
- [82] E. Jabason, M. O. Ahmad, and M. N. S. Swamy, "Missing Structural and Clinical Features Imputation for Semi-supervised Alzheimer's Disease Classification using Stacked Sparse Autoencoder," in *2018 IEEE Biomedical Circuits and Systems Conference (BioCAS)*, 17-19 Oct. 2018 2018, pp. 1-4, doi: 10.1109/BIOCAS.2018.8584844.
- [83] M. Mirza and S. Osindero, "Conditional Generative Adversarial Nets," *arXiv e-prints*, pp. arXiv:1411.1784-arXiv:1411.1784, 2014.
- [84] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks," *arXiv e-prints*, pp. arXiv:1703.10593-arXiv:1703.10593, 2017.
- [85] A. Dosovitskiy *et al.*, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *arXiv e-prints*, pp. arXiv:2010.11929-arXiv:2010.11929, 2020.
- [86] Z. Liu *et al.*, "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," *arXiv e-prints*, pp. arXiv:2103.14030-arXiv:2103.14030, 2021.
- [87] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," *arXiv e-prints*, pp. arXiv:1810.04805-arXiv:1810.04805, 2018.
- [88] J. Lee *et al.*, "BioBERT: a pre-trained biomedical language representation model for biomedical text mining," *Bioinformatics*, <https://doi.org/10.1093/bioinformatics/btz682> vol. 36, no. 4, pp. 1234-1240, 2019. [Online]. Available: <https://doi.org/10.1093/bioinformatics/btz682>.

- [89] L. Rasmy, Y. Xiang, Z. Xie, C. Tao, and D. Zhi, "Med-BERT: pretrained contextualized embeddings on large-scale structured electronic health records for disease prediction," *npj Digital Medicine*, vol. 4, no. 1, p. 86, 2021/05/20 2021, doi: 10.1038/s41746-021-00455-y.
- [90] Y. Luo, X. Cai, Y. Zhang, J. Xu, and Y. xiaojie, "Multivariate Time Series Imputation with Generative Adversarial Networks," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., 2018, vol. 31: Curran Associates, Inc. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/file/96b9bff013acedfb1d140579e2fbeb63-Paper.pdf>. [Online]. Available: <https://proceedings.neurips.cc/paper/2018/file/96b9bff013acedfb1d140579e2fbeb63-Paper.pdf>
- [91] C. H. Lee and H.-J. Yoon, "Medical big data: promise and challenges," (in eng), *Kidney Res Clin Pract*, vol. 36, no. 1, pp. 3-11, 2017, doi: 10.23876/j.krcp.2017.36.1.3.
- [92] K. Huang, C. Xiao, L. M. Glass, C. W. Critchlow, G. Gibson, and J. Sun, "Machine learning applications for therapeutic tasks with genomics data," *Patterns*, vol. 2, no. 10, p. 100328, 2021/10/08/ 2021, doi: <https://doi.org/10.1016/j.patter.2021.100328>.
- [93] E. R. David and L. M. James, "Learning Internal Representations by Error Propagation," in *Parallel Distributed Processing: Explorations in the Microstructure of Cognition: Foundations*: MIT Press, 1987, pp. 318-362.
- [94] M. A. Kramer, "Nonlinear principal component analysis using autoassociative neural networks," *AIChE Journal*, vol. 37, no. 2, pp. 233-243, 1991, doi: <https://doi.org/10.1002/aic.690370209>.
- [95] I. Goodfellow *et al.*, "Generative Adversarial Nets," in *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Q. Weinberger, Eds., 2014, vol. 27: Curran Associates, Inc. [Online]. Available: <https://proceedings.neurips.cc/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf>. [Online]. Available: <https://proceedings.neurips.cc/paper/2014/file/5ca3e9b122f61f8f06494c97b1afccf3-Paper.pdf>
- [96] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.

Table 1 Definitions of models and corresponding imputation mechanisms

Model	Definition	Imputation Mechanism ¹
Multi-layer perceptron (MLP)	It consists of at least three layers of neurons (input, hidden, and output), each of which is fed a non-linear activation function to capture the patterns.[92]	 <p>The diagram illustrates the MLP imputation mechanism. It starts with 'Original data' represented as a grid. This data is used for 'Training & Validation' and 'Testing'. The model structure shows an input layer, several 'Hidden layers', and an output layer. The 'Prediction' phase uses the trained model to generate 'Imputed data', where missing values are filled in the original data grid.</p>
Recurrent neural network (RNN) ²	It uses the hidden states (memory) to process sequences of inputs. At each step t , the input includes the current observation (x_t) and the previous hidden state (h_{t-1}).[93]	 <p>(a) shows a sequence of hidden states $h_{t-1}, h_t, h_{t+1}, \dots, h_T$ and observations $x_{t-1}, x_t, x_{t+1}, \dots, x_T$. Forward predictions \hat{x}_t^f and forward imputations \hat{x}_t^f are shown. (b) shows a detailed view of a hidden state h_{t-1} with gates and activation functions (σ sigmoid, \tanh tanh) and operations (+ plus, \otimes Hadamard product).</p> <p>Legend: σ sigmoid \oplus plus \tanh tanh \otimes Hadamard product</p>
Autoencoder (AE) ³	It is a framework ⁴ that co-trains two modules: an encoder that maps the input data into a latent embedding, and a decoder that reconstructs the input from the latent embedding.[94]	 <p>The diagram illustrates the Autoencoder imputation mechanism. 'Original data' is processed by an 'Encoder' to produce a 'Latent Space'. The 'Latent Space' is then processed by a 'Decoder' to produce 'Reconstructed data'. The reconstructed data is used to generate 'Imputed data', where missing values are filled in the original data grid. A legend indicates that hatched areas represent 'Random values'.</p>



¹Imputation mechanisms: MLP models can be trained on the complete observations to predict the missing values; RNN models can predict the missing values based on the previous hidden state (forward imputation[45]); Autoencoder can maintain the whole data structure in a good manner and reveal the missing values in its output; GAN can use the generator to capture the data distribution, impute the missing values with the generated data, and apply the discriminator to decide the rightness of the imputation with the assistance of domain knowledge, if applicable. The adversarial process allows for precise data distribution capturing. Based on these general ideas, applications and variants are discussed in subsections 3.1-3.6.

²Long short-term memory (LSTM) and Gated recurrent unit (GRU) are two main branches of RNN. Compared with vanilla RNN, they have an additional mechanism of “gates” to control the contribution of “memory”, i.e., sequence-dependencies, as shown in (b). GRU with two gates is simpler than LSTM with three gates, but performs similarly in many scenarios.

³Denoise autoencoder (DAE) and variational autoencoder (VAE) hold the same fundamental structure as vanilla autoencoder. DAE receives corrupted data points as input and is trained to predict the uncorrupted data points as its output[96]. Considering missingness as one of the forms of corruption, DAE can be more robust to missing values than vanilla AE. Variational Autoencoder (VAE) utilizes the technique of variational inference in statistics, which introduces probabilistic modeling in latent space to better approximate the true data[61].

⁴Framework refers to the fact that the modules (encoder/decoder and generator/discriminator) which respectively shaped the structures of AE and GAN, can embed with various models based on the data input, for example, convolutional neural network (CNN) to tackle images.

Figure 1 Preferred Reporting Items for Systematic Reviews (PRISMA) flow diagram

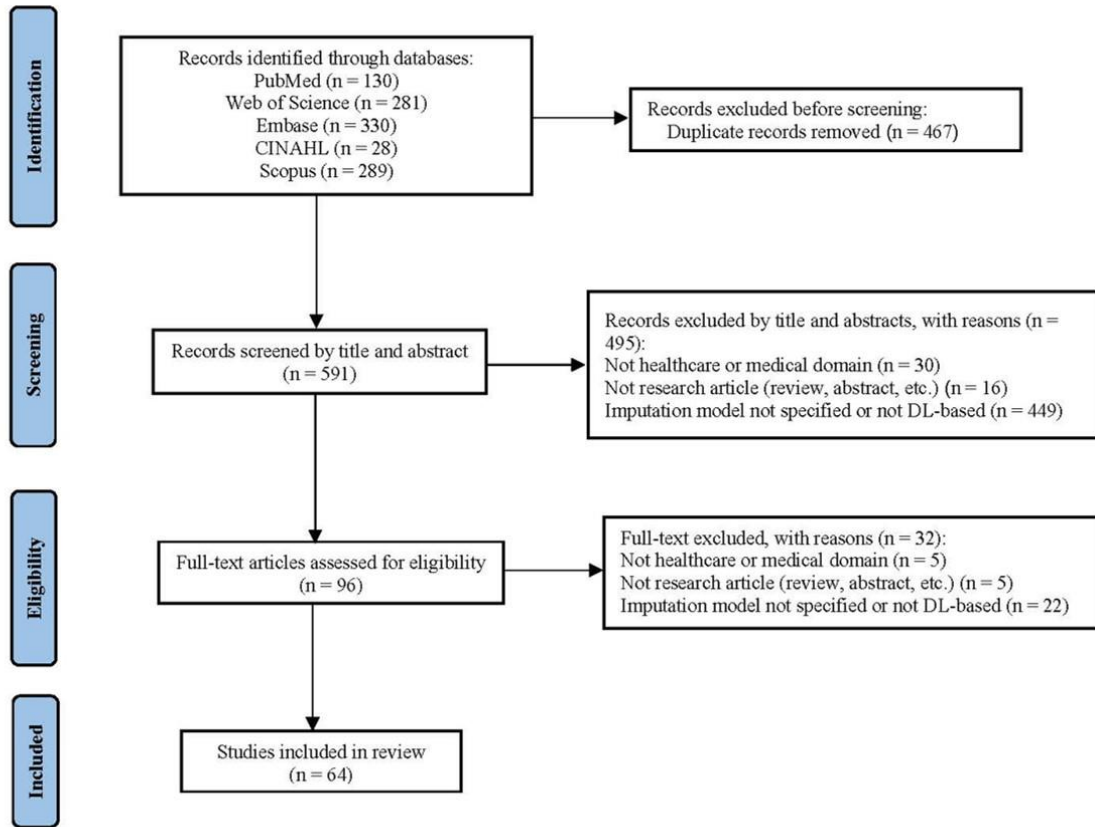
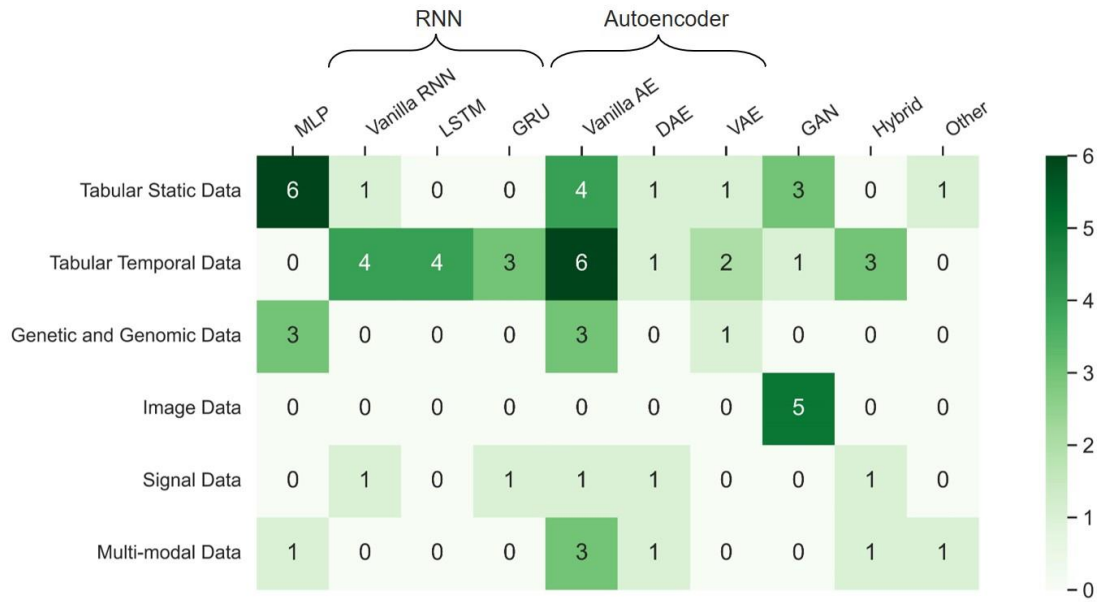
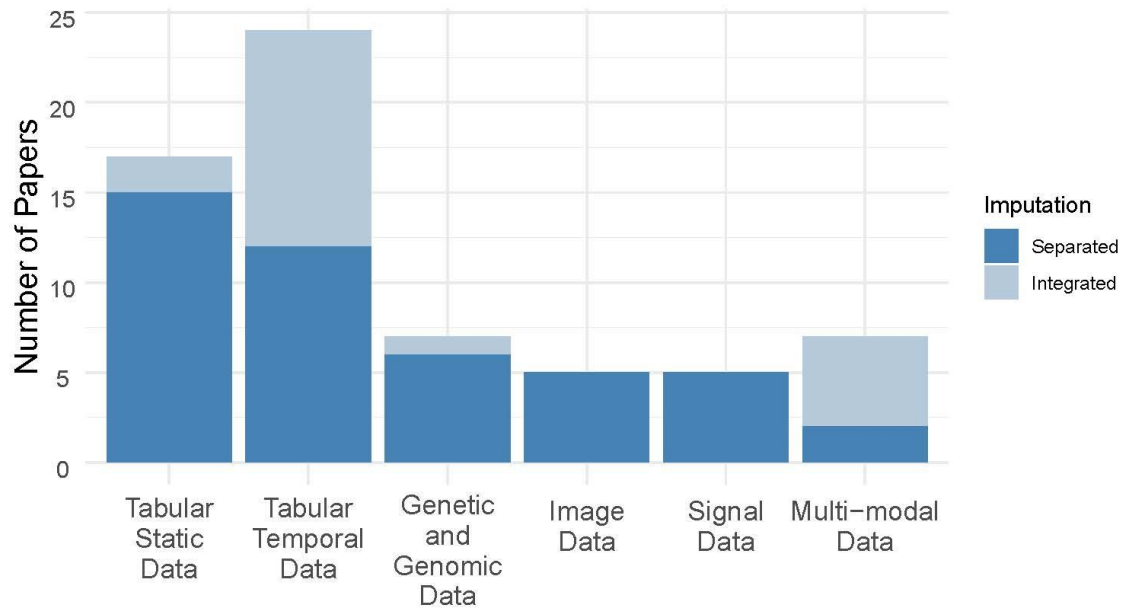


Figure 2 Evidence map between "backbones" (main architectures) of model and data type



"Backbones" are classified into ten categories: MLP (multi-layer perceptron), RNN (recurrent neural network), LSTM (Long short-term memory), GRU (gated recurrent unit), AE (autoencoder), DAE (denoising autoencoder), VAE (variational autoencoder), GAN (generative adversarial network) and Other, which includes less frequently used models such as SOM (self-organizing map). Data types are categorized into seven categories: tabular static, tabular temporal, genetic and genomic, image, signal, and multi-modal data. The numbers are non-exclusive.

Figure 3 The distribution of imputation frameworks (integrated or separated) by data type



Supplementary

eTable 1. Search Strategies

Database	Search Strategy
MEDLINE	("machine learning" OR "deep learning" OR "artificial intelligence" OR "data mining") AND ("missing data" OR "missing value" OR "missing values") AND ("imputation") AND ("medical" OR "clinical" OR "health" OR "healthcare" OR "medicine")
Web of Science	ALL=("machine learning" OR "deep learning" OR "artificial intelligence" OR "data mining") AND ALL=("missing data" OR "missing value" OR "missing values") AND ALL=("imputation") AND ALL=("medical" OR "clinical" OR "health" OR "healthcare" OR "medicine")
Embase	("machine learning" OR "deep learning" OR "artificial intelligence" OR "data mining") AND ("missing data" OR "missing value" OR "missing values") AND ("imputation") AND ("medical" OR "clinical" OR "health" OR "healthcare" OR "medicine")
CINAHL	("machine learning" OR "deep learning" OR "artificial intelligence" OR "data mining") AND ("missing data" OR "missing value" OR "missing values") AND ("imputation") AND ("medical" OR "clinical" OR "health" OR "healthcare" OR "medicine")
Scopus	("machine learning" OR "deep learning" OR "artificial intelligence" OR "data mining") AND ("missing data" OR "missing value" OR "missing values") AND ("imputation") AND ("medical" OR "clinical" OR "health" OR "healthcare" OR "medicine")