

Contrastive Learning-based Imputation-Prediction Networks for In-hospital Mortality Risk Modeling using EHRs

Yuxi Liu✉¹, Zhenhao Zhang², Shaowen Qin¹, Flora D. Salim³, and Antonio Jimeno Yepes⁴

¹ College of Science and Engineering, Flinders University, Tonsley, SA 5042, Australia {liu1356, shaowen.qin}@flinders.edu.au

² College of Life Sciences, Northwest A&F University, Yangling, Shaanxi 712100, China {zhangzhenhao}@nwfufu.edu.cn

³ School of Computer Science and Engineering, UNSW, Sydney, NSW 2052, Australia {flora.salim}@unsw.edu.au

⁴ School of Computing Technologies, RMIT University, Melbourne, VIC 3001, Australia {antonio.jose.jimeno.yepes}@rmit.edu.au

Abstract. Predicting the risk of in-hospital mortality from electronic health records (EHRs) has received considerable attention. Such predictions will provide early warning of a patient’s health condition to healthcare professionals so that timely interventions can be taken. This prediction task is challenging since EHR data are intrinsically irregular, with not only many missing values but also varying time intervals between medical records. Existing approaches focus on exploiting the variable correlations in patient medical records to impute missing values and establishing time-decay mechanisms to deal with such irregularity. This paper presents a novel contrastive learning-based imputation-prediction network for predicting in-hospital mortality risks using EHR data. Our approach introduces graph analysis-based patient stratification modeling in the imputation process to group similar patients. This allows information of similar patients only to be used, in addition to personal contextual information, for missing value imputation. Moreover, our approach can integrate contrastive learning into the proposed network architecture to enhance patient representation learning and predictive performance on the classification task. Experiments on two real-world EHR datasets show that our approach outperforms the state-of-the-art approaches in both imputation and prediction tasks.

Keywords: data imputation · in-hospital mortality · contrastive learning.

1 Introduction

The broad adoption of digital healthcare systems produces a large amount of electronic health records (EHRs) data, providing us the possibility to develop

predictive models and tools using machine learning techniques that would enable healthcare professionals to make better decisions and improve healthcare outcomes. One of the EHR-based risk prediction tasks is to predict the mortality risk of patients based on their historical EHR data [8,29]. The predicted mortality risks can be used to provide early warnings when a patient’s health condition is about to deteriorate so that more proactive interventions can be taken.

However, due to a high degree of irregularity in the raw EHR data, it is challenging to directly apply traditional machine learning techniques to perform predictive modeling. We take the medical records of two anonymous patients from the publicly available MIMIC-III database and present these in Figure 1 as an example. Figure 1 clearly indicates the irregularity problem, including many missing values and varying time intervals between medical records.

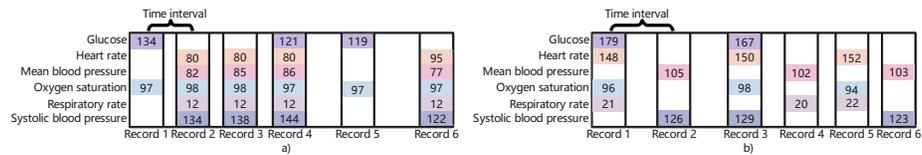


Fig. 1. Illustration of medical records of patients A and B.

Most studies have focused on exploiting variable correlations in patient medical records to impute missing values and establishing time-decay mechanisms to take into account the effect of varying time intervals between records [1, 2, 17, 18, 23–25, 31]. After obtaining the complete data matrices from the imputation task, the complete data matrices are used as input for downstream healthcare prediction tasks [1, 2, 13, 17, 18, 22, 23, 27, 30, 31, 35]. Although these studies have achieved satisfactory imputation performance, consideration of using the information of similar patients on the imputation task, which might lead to improved imputation performance, has not yet been fully experimented. Furthermore, with imputation data, high-quality representation must be applied, as the imputation data may affect the performance of downstream healthcare prediction tasks.

Patient stratification refers to the method of dividing a patient population into subgroups based on specific disease characteristics and symptom severity. Patients in the same subgroup generally had more similar health trajectories. Therefore, we propose to impute missing values in patient data using information from the subgroup of similar patients rather than the entire patient population.

In this paper, we propose a novel contrastive learning-based imputation-prediction network with the aim of improving in-hospital mortality prediction performance using EHR data. Missing value imputation for EHR data is done by exploiting similar patient information as well as patients’ personal contextual information. Similar patients are generated from patient similarity calculation during stratification modeling and analysis of patient graphs.

Contrastive learning has been proven to be an important machine learning technique in the computer vision community [12]. In contrastive learning, representations are learned by comparing input samples. The comparisons are made on the similarity between positive pairs or dissimilarity between negative pairs. The main goal is to learn an embedding space where similar samples are put closer to each other while dissimilar samples are pushed farther apart. Contrastive learning can be applied in both supervised [10, 33, 39] and unsupervised [14, 15, 26] settings.

Motivated by the recent developments in contrastive representation learning [34, 36, 38], we integrate contrastive learning into the proposed network architecture to perform imputation and prediction tasks. The benefit of incorporating contrastive learning into the imputation task is that such an approach can enhance patient representation learning by keeping patients of the same stratification together and pushing away patients from different stratifications. This would lead to enhanced imputation performance. The benefit of incorporating contrastive learning into the prediction task is improved predictive performance of the binary classification problem (i.e., the risk of death and no death), which is achieved by keeping the instances of a positive class closer and pushing away instances from a negative class.

Our major contributions are as follows:

- To the best of our knowledge, this is the first attempt to consider patient similarity via stratification of EHR data on the imputation task.
- We propose a novel imputation-prediction approach to perform imputation and prediction simultaneously with EHR data.
- We successfully integrate contrastive learning into the proposed network architecture to improve imputation and prediction performance.
- Extensive experiments conducted on two real-world EHR datasets show that our approach outperforms all baseline approaches in imputation and prediction tasks.

2 Related Work

There has been an increased interest in EHR-based health risk predictions [5, 16, 19–21]. It has been recognized that EHR data often contains many missing values due to patient conditions and treatment decisions [31]. Existing research addresses this challenge by imputing missing data and feeding them into the supervised algorithms as auxiliary information [7]. GRU-D [2] represents such an example. The GRU-D is built upon the Gated Recurrent Unit [4]. GRU-D proposes to impute missing values by decaying the contributions of previous observation values toward the overall mean over time. Similarly, BRITS [1] incorporates a bidirectional recurrent neural network (RNN) to impute missing values. Since the incorporated bidirectional RNN learns EHR data in both forward and backward directions, the accumulated loss is introduced to train the model.

Another line of related work is based on the generative adversarial network (GAN) architecture, which aims at treating the problem of missing data imputation as data generation. The intuitions behind GAN can be seen as making a generator and a discriminator against each other [6]. The generator generates fake samples from random ‘noise’ vectors, and the discriminator distinguishes the generator’s fake samples from actual samples. Examples of research into GAN-based imputation methods include GRUI-GAN [17], E²GAN [18], E²GAN-RF [40], and STING [25]. These studies take the vector of actual samples, which has many missing values, use a generator to generate the corresponding imputed values and distinguish the generated imputed values from real values using a discriminator.

Several studies have evaluated the effectiveness of applying transformer-based imputation methods to EHR data. Examples of representative studies include MTSIT [37] and MIAM [13]. The MTSIT is built with an autoencoder architecture to perform missing value imputation in an unsupervised manner. The autoencoder architecture used in MTSIT includes the Transformer encoder [32] and a linear decoder, which are implemented with a joint reconstruction and imputation approach. The MIAM is built upon the self-attention mechanism [32]. Given EHR data, MIAM imputes the missing values by extracting the relationship among the observed values, missingness indicators (0 for missing and 1 for not missing), and the time interval between consecutive observations.

3 Method

3.1 Network Architecture

The architecture of the proposed network is shown in Figure 2.

Data Representation We represent a multivariate time series X with up to N variables of length T as a set of observed triplets, i.e., $X = \{(f_i, v_i, t_i)\}_{i=1}^N$. An observed triplet is represented as a (f, v, t) , where $f \in F$ is the variable/feature, $v \in \mathbb{R}^T$ is the observed value, and $t \in \mathbb{R}^T$ is the time. We incorporate a masking vector m_i to represent missing values in v_i as:

$$m_{i,t} = \begin{cases} 1, & \text{if } v_{i,t} \text{ is observed} \\ 0, & \text{otherwise} \end{cases} \quad (1)$$

Let $\delta \in \mathbb{R}^{N \times T}$, $\delta^{(l)} \in \mathbb{R}^{N \times T}$, and $\delta^{(n)} \in \mathbb{R}^{N \times T}$ denote three time interval matrices. δ_t is the time interval between the current time t and the last time $t - 1$. $\delta_{i,t}^{(l)}$ is the time interval between the current time t and the time where the i -th variable is observed the last time. $\delta_{i,t}^{(n)}$ is the time interval between the current time t and the time where the i -th variable is observed next time. $\delta_{i,t}^{(l)}$ and $\delta_{i,t}^{(n)}$ can be written as:

$$\delta_{i,t}^{(l)} = \begin{cases} \delta_{i,t}, & \text{if } m_{i,t-1} = 1 \\ \delta_{i,t} + \delta_{i,t-1}^{(l)}, & \text{otherwise} \end{cases} \quad (2)$$

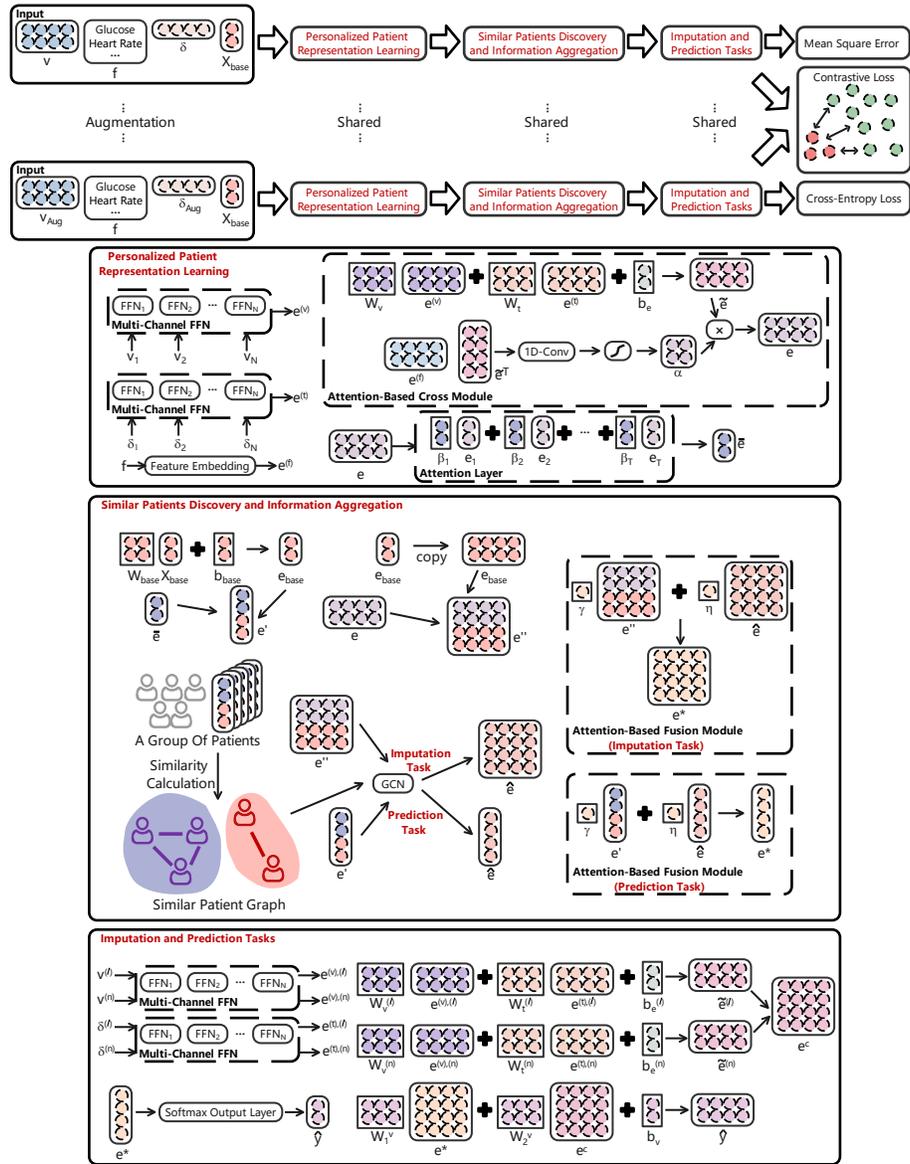


Fig. 2. Schematic description of the proposed network.

$$\delta_{i,t}^{(n)} = \begin{cases} \delta_{i,t+1}, & \text{if } m_{i,t+1} = 1 \\ \delta_{i,t+1} + \delta_{i,t+1}^{(n)}, & \text{otherwise} \end{cases} \quad (3)$$

Let $v^{(l)}$ and $v^{(n)}$ denote two neighboring value matrices, the observed values of the last time and next time. $v^{(l)}$ and $v^{(n)}$ can be written as:

$$v_{i,t}^{(l)} = \begin{cases} v_{i,t-1}, & \text{if } m_{i,t-1} = 1 \\ v_{i,t-1}^{(l)}, & \text{otherwise} \end{cases} \quad (4)$$

$$v_{i,t}^{(n)} = \begin{cases} v_{i,t+1}, & \text{if } m_{i,t+1} = 1 \\ v_{i,t+1}^{(n)}, & \text{otherwise} \end{cases} \quad (5)$$

where $v_{i,t}^{(l)}$ and $v_{i,t}^{(n)}$ are the values of the i -th variable of $v_t^{(l)}$ and $v_t^{(n)}$.

Let $D = \{(X_p, y_p)\}_{p=1}^P$ denote the EHR dataset with up to P labeled samples. The p -th sample contains a multivariate time series X_p consisting of the physiological variables, and a binary label of in-hospital mortality $y_p \in \{0, 1\}$. Let $X_{base} \in \mathbb{R}^g$ denote the patient-specific characteristics (i.e., age, sex, ethnicity, admission diagnosis) with up to g dimension.

Personalized Patient Representation Learning Given an input multivariate time series/a single patient data $X = \{(f_i, v_i, t_i)\}_{i=1}^N$, the embedding for the i -th triplet $e_i \in \mathbb{R}^d$ is generated by aggregating the feature embedding $e_i^{(f)} \in \mathbb{R}^d$, the value embedding $e_i^{(v)} \in \mathbb{R}^{d \times T}$, and the time interval embedding $e_i^{(t)} \in \mathbb{R}^{d \times T}$. The feature embedding is similar to the word embedding, which allows features with similar meanings to have a similar representation. Particularly, the value embedding and time interval embedding are obtained by separately implementing a multi-channel feed-forward neural network (FFN) as:

$$\begin{aligned} e_{i,1}^{(v)}, \dots, e_{i,T}^{(v)} &= FFN_i^{(v)}(v_{i,1}, \dots, v_{i,T}), \\ e_{i,1}^{(t)}, \dots, e_{i,T}^{(t)} &= FFN_i^{(t)}(\delta_{i,1}, \dots, \delta_{i,T}). \end{aligned} \quad (6)$$

Through the processes above, we are able to obtain $e^{(f)} \in \mathbb{R}^{Nd}$, $e^{(v)} \in \mathbb{R}^{Nd \times T}$, and $e^{(t)} \in \mathbb{R}^{Nd \times T}$, which are fed into the attention-based cross module to generate an overall representation. Note that $e^{(f)} \in \mathbb{R}^{Nd}$ is expanded into $e^{(f)} \in \mathbb{R}^{Nd \times T}$. Specifically, we design the attention-based cross module to generate a cross-attention matrix as:

$$\begin{aligned} \tilde{e} &= W_v \cdot e^{(v)} + W_t \cdot e^{(t)} + b_e, \\ E &= ScaledDot(e^{(f)}, \tilde{e}) = \frac{e^{(f)} \cdot \tilde{e}^\top}{\sqrt{d}}, \end{aligned} \quad (7)$$

where $E \in \mathbb{R}^{Nd \times Nd}$ is the cross-attention matrix that corresponds to the scaled-dot similarity. We then apply a 1D convolutional layer to the cross-attention matrix E as:

$$\alpha = \text{Softmax}(\text{Conv}(E)), \quad (8)$$

where Conv is the 1D convolutional layer and α is the cross-attention score matrix. We integrate α and \tilde{e} into a weighted representation e as:

$$e = \alpha \odot \tilde{e}. \quad (9)$$

Given a batch of patients, the embedding for them can be written as:

$$e = [e_1, e_2, \dots, e_B] \in \mathbb{R}^{B \times Nd \times T}, \quad (10)$$

where B is the batch size. Since e still takes the form of sequence data, we design an attention layer to generate a series of attention weights $(\beta_1, \beta_2, \dots, \beta_T)$ and reweight these weights to produce an overall feature representation as:

$$\begin{aligned} \beta &= \text{Softmax}(e \cdot W_e + b_e), \\ \bar{e} &= \sum_{t=1}^T \beta_t \odot e_t, \end{aligned} \quad (11)$$

where $\bar{e} \in \mathbb{R}^{B \times Nd}$ is the new generated patient representation.

Similar Patients Discovery and Information Aggregation Before conducting patient similarity calculation, we encode $X_{base} \in \mathbb{R}^g$ as $e_{base} \in \mathbb{R}^{d_g}$ and concatenate e_{base} with \bar{e} as:

$$\begin{aligned} e_{base} &= W_{base} \cdot X_{base} + b_{base}, \\ e' &= \text{Concat}(\bar{e}, e_{base}), \end{aligned} \quad (12)$$

where Concat is the concatenation operation.

For the batch of patient representations, the pairwise similarities that correspond to any two patient representations can be calculated as:

$$\Lambda = \text{sim}(e', e') = \frac{e' \cdot e'}{(Nd + d_g)^2}, \quad (13)$$

where $\text{sim}(\cdot)$ is the measure of cosine similarity and $\Lambda \in \mathbb{R}^{B \times B}$ is the patient similarity matrix.

Moreover, we incorporate a learnable threshold φ into the patient similarity calculation to filter out similarities below the threshold. The similarity matrix can be rewritten as:

$$\Lambda' = \begin{cases} \Lambda, & \text{if } \Lambda > \varphi \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

We take into account the batch of patients' representations as a graph to aggregate the information from similar patients, where the similarity matrix A' is the graph adjacency matrix. We apply graph convolutional layers to enhance the representation learning as:

$$\begin{aligned}\hat{e} &= [\hat{e}_1, \hat{e}_2, \dots, \hat{e}_B]^\top = GCN(e', A') \\ &= ReLU(A' ReLU(A' \cdot e' W_1^e) \cdot W_2^e),\end{aligned}\tag{15}$$

where \hat{e} is the aggregated auxiliary information from similar patients. A note of caution is due here since we ignore the bias term. We replace e' in Eq. (15) with e'' for the imputation task. By doing so, the output of graph convolutional layers can take the form of sequence data. Particularly, e'' is obtained by concatenating e and e_{base} , where $e_{base} \in \mathbb{R}^{d_g}$ is expanded into $e_{base} \in \mathbb{R}^{d_g \times T}$.

Through the processes above, we are able to generate e'/e'' and \hat{e} representations for the batch of patients. The e'/e'' refers to the patient themselves. For an incomplete patient p (i.e., the patient data has many missing values), we generate the missing value representations with \hat{e} . For a complete patient, we augment e'/e'' with \hat{e} to enhance the representation learning.

We design an attention-based fusion module to refine both e'/e'' (the two representations used in prediction and imputation tasks) and \hat{e} . Since imputation and prediction tasks involve the same process of modeling, we take the prediction task as an example. The two weights $\gamma \in \mathbb{R}^B$ and $\eta \in \mathbb{R}^B$ are incorporated to determine the importance of e' and \hat{e} , obtained by implementing fully connected layers as:

$$\begin{aligned}\gamma &= Sigmoid(e' \cdot W_\gamma + b_\gamma), \\ \eta &= Sigmoid(\hat{e} \cdot W_\eta + b_\eta).\end{aligned}\tag{16}$$

A note of caution is due here since we keep the sum of γ and η must be 1, i.e., $\gamma + \eta = 1$. We achieve this constraint by combining $\gamma = \frac{\gamma}{\gamma + \eta}$ and $\eta = 1 - \gamma$. The final representation e^* is obtained by calculating $\gamma \cdot e' + \eta \cdot \hat{e}$.

Contrastive Learning We integrate contrastive learning into the proposed network architecture to perform imputation and prediction tasks. For the prediction task, we augment the standard cross-entropy loss with the supervised contrastive loss [10]. We treat the patient representations with the same label as the positive pairs and the patient representations with different labels as the negative pairs. For the imputation task, we augment the standard mean squared error loss with the unsupervised contrastive loss [3]. We treat a single patient representation and its augmented representations as positive pairs and the other patient representations within a batch and their augmented representations as negative pairs. The formula can be written as:

$$\begin{aligned}\mathcal{L}_{SC} &= - \sum_{i=1}^B \frac{1}{B_{y_i}} \log \frac{\sum_{j=1}^B \mathbb{1}_{[y_i=y_j]} \exp(sim(e_i^*, e_j^*)/\tau)}{\sum_{k=1}^B \mathbb{1}_{[k \neq i]} \exp(sim(e_i^*, e_k^*)/\tau)}, \\ \mathcal{L}_{UC} &= - \log \frac{\exp(sim(e_i^*, e_j^*)/\tau)}{\sum_{k=1}^{2B} \mathbb{1}_{[k \neq i]} \exp(sim(e_i^*, e_k^*)/\tau)},\end{aligned}\tag{17}$$

where B represents the batch size; $\mathbb{1}_{[\cdot]}$ represents an indicator function; $sim(\cdot)$ represents the cosine similarity measure; τ represents a hyper-parameter that is used to control the strength of penalties on negative pairs; B_{y_i} is the number of samples with the same label in each batch.

Imputation and Prediction Tasks For the prediction task, we feed e^* into a softmax output layer to obtain the predicted \hat{y} as:

$$\hat{y} = \text{Softmax}(W_y \cdot e^* + b_y). \quad (18)$$

The objective loss is the summation of cross-entropy loss and the supervised contrastive loss with a scaling parameter λ to control the contribution of each loss as:

$$\begin{aligned} \mathcal{L}_{CE} &= -\frac{1}{P} \sum_{p=1}^P (y_p^\top \cdot \log(\hat{y}_p) + (1 - y_p)^\top \cdot \log(1 - \hat{y}_p)), \\ \mathcal{L} &= \lambda \cdot \mathcal{L}_{CE} + (1 - \lambda) \cdot \mathcal{L}_{SC}. \end{aligned} \quad (19)$$

For the imputation task, we take the neighboring observed values (of each patient) as inputs to incorporate patient-specific contextual information. The process of embedding used by $v^{(l)}$ and $v^{(n)}$ can be written as:

$$\begin{aligned} e_i^{(v),(l)} &= FFN_i^{(v),(l)}(v_i^{(l)}), e_i^{(t),(l)} = FFN_i^{(t),(l)}(\delta_i^{(l)}), \\ e_i^{(v),(n)} &= FFN_i^{(v),(n)}(v_i^{(n)}), e_i^{(t),(n)} = FFN_i^{(t),(n)}(\delta_i^{(n)}), \\ \tilde{e}^{(l)} &= W_v^{(l)} \cdot e^{(v),(l)} + W_t^{(l)} \cdot e^{(t),(l)} + b_e^{(l)}, \\ \tilde{e}^{(n)} &= W_n^{(v)} \cdot e^{(v),(n)} + W_t^{(n)} \cdot e^{(t),(n)} + b_e^{(n)}, \\ e^c &= \text{Concate}(\tilde{e}^{(l)}, \tilde{e}^{(n)}), \end{aligned} \quad (20)$$

where $\tilde{e}^{(l)}$ and $\tilde{e}^{(n)}$ are the representations of $v^{(l)}$ and $v^{(n)}$ after embedding. The embedding matrix e^c is obtained by concatenating $\tilde{e}^{(l)}$ and $\tilde{e}^{(n)}$.

Given the final representation e^* and the embedding matrix e^c , we use a fully connected layer to impute missing values as:

$$\hat{v} = e^* \cdot W_1^v + e^c \cdot W_2^v + b_v. \quad (21)$$

The objective loss is the summation of the mean square error and the unsupervised contrastive loss with a scaling parameter λ to control the contribution of each loss as:

$$\begin{aligned} \mathcal{L}_{MSE} &= \frac{1}{P} \sum_{p=1}^P (m_p \odot v_p - m_p \odot \hat{v}_p)^2, \\ \mathcal{L} &= \lambda \cdot \mathcal{L}_{MSE} + (1 - \lambda) \cdot \mathcal{L}_{UC}. \end{aligned} \quad (22)$$

4 Experiments

4.1 Datasets and Tasks

We validate our approach on the MIMIC-III⁵ and eICU⁶ datasets. We conduct clinical time series imputation and in-hospital mortality experiments based on the data from the first 24/48 hours after admission. Detailed information on both datasets can be found in the literature [9] and [28]. The source code of our approach and statistics of features are released at <https://github.com/liulab1356/CL-ImpPreNet>.

4.2 Baseline Approaches

We compare our approach with GRU-D [2], BRITS [1], GRUI-GAN [17], E²GAN [18], E²GAN-RF [40], STING [25], MTSIT [37], and MIAM [13] (see related work section). We feed the output of GRUI-GAN, E²GAN, E²GAN-RF, STING, and MTSIT into GRU to estimate in-hospital mortality risk probabilities. Moreover, the regression component used in BRITS is integrated into GRU-D and MIAM to obtain imputation accuracy.

Besides, two variants of our approach are as follows:

Ours_α: We do not perform graph analysis-based patient stratification modeling.

Ours_β: We omit the contrastive learning component.

All implementations of Ours_α and Ours_β can be found in the aforementioned Github repository.

4.3 Implementation Details and Evaluation Metrics

We implement all approaches with PyTorch 1.11.0 and conduct experiments on A40 GPU from NVIDIA with 48GB of memory. We randomly use 70%, 15%, and 15% of the dataset as training, validation, and testing sets. We train the proposed approach using an Adam optimizer [11] with a learning rate of 0.0023 and a mini-batch size of 256. For personalized patient representation learning, the dimension size d is 3. For similar patients discovery and information aggregation, the initial value of φ is 0.56, and the dimension size of W_1^e and W_2^e are 34 and 55. For contrastive learning, the value of τ is 0.07. The dropout method is applied to the final Softmax output layer for the prediction task, and the dropout rate is 0.1. For the imputation task, the dimension size of $W_v^{(l)}$, $W_t^{(l)}$, $W_v^{(n)}$, and $W_t^{(n)}$ are 28.

The performance of contrastive learning heavily relies on data augmentation. We augment the observed value v with random time shifts and reversion. For example, given the observed value $v = [v_1, v_2, \dots, v_T]$, we are able to obtain $v_{shift} = [v_{1+n}, v_{2+n}, \dots, v_{T+n}]$ and $v_{reverse} = [v_T, v_{T-1}, \dots, v_1]$ from random time shift and reversion, and n is the number of data points to shift.

⁵ <https://mimic.physionet.org>

⁶ <https://eicu-crd.mit.edu/>

We use the MAE and MRE scores between predicted and actual values as the evaluation metrics for imputation performance. We use the AUROC and AUPRC scores as the evaluation metrics for prediction performance. We report the mean and standard deviation of the evaluation metrics after repeating all the approaches ten times.

5 Experimental Results

Table 2 presents the experimental results of all approaches on imputation and prediction tasks from MIMIC-III and eICU datasets. Together these results suggest that our approach achieves the best performance in both imputation and prediction tasks. For example, for the clinical time series imputation of MIMIC-III (24 hours after ICU admission), the MAE and MRE of Ours are 0.3563 and 8.16%, smaller than 0.3988 and 38.44% achieved by the best baseline (i.e., MTSIT). For the in-hospital mortality prediction of MIMIC-III (24 hours after ICU admission), the AUROC and AUPRC of Ours are 0.8533 and 0.4752, larger than 0.8461 and 0.4513 achieved by the best baseline (i.e., GRU-D).

As Table 2 shows, the RNN-based approach (i.e., GRU-D and BRITS) outperforms the GAN-based approach (i.e., GRUI-GAN, E²GAN, E²GAN-RF, and STING) in the imputation task. From the prediction results of the MIMIC-III dataset, we can see that the transformer-based approaches (i.e., MTSIT and MIAM) resulted in lower values of AUROC and AUPRC. From the prediction results of the eICU dataset, no significant difference between the transformer-based approach and other approaches was evident.

Ours outperforms its variants Ours _{α} and Ours _{β} . This result confirms the effectiveness of the network construction with enhanced imputation and prediction performance.

6 Conclusion

This paper presents a novel contrastive learning-based imputation-prediction network to carry out in-hospital mortality prediction tasks using EHR data. This prediction makes timely warnings available to ICU health professionals so that early interventions for patients at risk could take place. The proposed approach explicitly considers patient similarity by stratification of EHR data and successfully integrates contrastive learning into the network architecture. We empirically show that the proposed approach outperforms all the baselines by conducting clinical time series imputation and in-hospital mortality prediction on the MIMIC-III and eICU datasets.

7 Acknowledgement

This research is partially funded by the ARC Centre of Excellence for Automated Decision-Making and Society (CE200100005) by the Australian Government through the Australian Research Council.

Table 1. Performance of our approaches with other baselines on clinical time series imputation and in-hospital mortality prediction.

| MIMIC-III/24 hours after ICU admission | Clinical time series imputation | | In-hospital mortality prediction | |
|--|---------------------------------|----------------------|----------------------------------|-----------------------|
| Metrics | MAE | MRE | AUROC | AUPRC |
| GRU-D | 1.3134(0.0509) | 87.33%(0.0341) | 0.8461(0.0051) | 0.4513(0.0124) |
| BRITS | 1.3211(0.0923) | 87.92%(0.0611) | 0.8432(0.0040) | 0.4193(0.0144) |
| GRUI-GAN | 1.6083(0.0043) | 107.20%(0.0029) | 0.8324(0.0077) | 0.4209(0.0280) |
| E ² GAN | 1.5885(0.0045) | 105.86%(0.0032) | 0.8377(0.0083) | 0.4295(0.0137) |
| E ² GAN-RF | 1.4362(0.0031) | 101.09%(0.0027) | 0.8430(0.0065) | 0.4328(0.0101) |
| STING | 1.5018(0.0082) | 102.53%(0.0047) | 0.8344(0.0126) | 0.4431(0.0158) |
| MTSIT | 0.3988(0.0671) | 38.44%(0.0647) | 0.8029(0.0117) | 0.4150(0.0165) |
| MIAM | 1.1391(0.0001) | 75.65%(0.0001) | 0.8140(0.0044) | 0.4162(0.0079) |
| Ours | 0.3563(0.0375) | 8.16%(0.0086) | 0.8533(0.0119) | 0.4752(0.0223) |
| Ours _α | 0.3833(0.0389) | 8.78%(0.0089) | 0.8398(0.0064) | 0.4555(0.0139) |
| Ours _β | 0.4125(0.0319) | 8.95%(0.0077) | 0.8417(0.0059) | 0.4489(0.0182) |
| eICU/24 hours after eICU admission | Clinical time series imputation | | In-hospital mortality prediction | |
| Metrics | MAE | MRE | AUROC | AUPRC |
| GRU-D | 3.9791(0.2008) | 52.11%(0.0262) | 0.7455(0.0107) | 0.3178(0.0190) |
| BRITS | 3.6879(0.3782) | 48.30%(0.0726) | 0.7139(0.0101) | 0.2511(0.0111) |
| GRUI-GAN | 9.1031(0.0130) | 119.29%(0.0016) | 0.7298(0.0094) | 0.3013(0.0141) |
| E ² GAN | 7.5746(0.0141) | 99.20%(0.0018) | 0.7317(0.0155) | 0.2973(0.0253) |
| E ² GAN-RF | 6.7108(0.0127) | 90.38%(0.0015) | 0.7402(0.0131) | 0.3045(0.0227) |
| STING | 7.1447(0.0651) | 93.56%(0.0083) | 0.7197(0.0154) | 0.2873(0.0182) |
| MTSIT | 1.6192(0.1064) | 21.20%(0.0139) | 0.7215(0.0071) | 0.2992(0.0115) |
| MIAM | 1.1726(0.3103) | 15.35%(0.0406) | 0.7262(0.0179) | 0.2659(0.0148) |
| Ours | 0.5365(0.0612) | 7.02%(0.0079) | 0.7626(0.0117) | 0.3388(0.0211) |
| Ours _α | 0.6792(0.0716) | 8.89%(0.0093) | 0.7501(0.0143) | 0.3325(0.0151) |
| Ours _β | 0.5923(0.0514) | 7.75%(0.0067) | 0.7533(0.0104) | 0.3303(0.0175) |
| MIMIC-III/48 hours after ICU admission | Clinical time series imputation | | In-hospital mortality prediction | |
| Metrics | MAE | MRE | AUROC | AUPRC |
| GRU-D | 1.4535(0.0806) | 86.47%(0.0482) | 0.8746(0.0026) | 0.5143(0.0077) |
| BRITS | 1.3802(0.1295) | 82.21%(0.0768) | 0.8564(0.0040) | 0.4445(0.0189) |
| GRUI-GAN | 1.7523(0.0030) | 104.50%(0.0018) | 0.8681(0.0077) | 0.5123(0.0166) |
| E ² GAN | 1.7436(0.0036) | 103.98%(0.0022) | 0.8705(0.0043) | 0.5091(0.0120) |
| E ² GAN-RF | 1.6122(0.0027) | 102.34%(0.0017) | 0.8736(0.0031) | 0.5186(0.0095) |
| STING | 1.6831(0.0068) | 100.46%(0.0035) | 0.8668(0.0123) | 0.5232(0.0236) |
| MTSIT | 0.4503(0.0465) | 30.42%(0.0314) | 0.8171(0.0114) | 0.4308(0.0189) |
| MIAM | 1.3158(0.0003) | 78.20%(0.0002) | 0.8327(0.0024) | 0.4460(0.0061) |
| Ours | 0.4396(0.0588) | 6.23%(0.0073) | 0.8831(0.0149) | 0.5328(0.0347) |
| Ours _α | 0.7096(0.0532) | 8.85%(0.0066) | 0.8671(0.0093) | 0.5161(0.0151) |
| Ours _β | 0.5786(0.0429) | 7.47%(0.0056) | 0.8709(0.0073) | 0.5114(0.0176) |
| eICU/48 hours after eICU admission | Clinical time series imputation | | In-hospital mortality prediction | |
| Metrics | MAE | MRE | AUROC | AUPRC |
| GRU-D | 5.8071(0.2132) | 44.53%(0.0164) | 0.7767(0.0141) | 0.3210(0.0182) |
| BRITS | 5.5546(0.5497) | 42.59%(0.0421) | 0.7285(0.0114) | 0.2510(0.0097) |
| GRUI-GAN | 14.0750(0.0301) | 107.96%(0.0021) | 0.7531(0.0167) | 0.2897(0.0201) |
| E ² GAN | 12.9694(0.0195) | 99.47%(0.0015) | 0.7605(0.0063) | 0.3014(0.0137) |
| E ² GAN-RF | 11.8138(0.0161) | 91.52%(0.0011) | 0.7763(0.0057) | 0.3101(0.0125) |
| STING | 12.0962(0.0806) | 92.79%(0.0062) | 0.7453(0.0182) | 0.2805(0.0190) |
| MTSIT | 2.8150(0.2105) | 21.58%(0.0161) | 0.7418(0.0091) | 0.3078(0.0120) |
| MIAM | 2.1146(0.4012) | 16.23%(0.0414) | 0.7574(0.0127) | 0.2776(0.0105) |
| Ours | 0.9412(0.0930) | 7.21%(0.0071) | 0.7907(0.0123) | 0.3417(0.0217) |
| Ours _α | 1.1099(0.1064) | 8.51%(0.0081) | 0.7732(0.0100) | 0.3311(0.0265) |
| Ours _β | 0.9930(0.0817) | 7.61%(0.0062) | 0.7790(0.0117) | 0.3335(0.0178) |

References

1. Cao, W., Wang, D., Li, J., Zhou, H., Li, L., Li, Y.: Brits: Bidirectional recurrent imputation for time series. *Advances in neural information processing systems* **31** (2018)
2. Che, Z., Purushotham, S., Cho, K., Sontag, D., Liu, Y.: Recurrent neural networks for multivariate time series with missing values. *Scientific reports* **8**(1), 1–12 (2018)
3. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: *International conference on machine learning*. pp. 1597–1607. PMLR (2020)
4. Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., Bengio, Y.: Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078* (2014)
5. Cui, S., Wang, J., Gui, X., Wang, T., Ma, F.: Automed: Automated medical risk predictive modeling on electronic health records. In: *2022 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. pp. 948–953. IEEE (2022)
6. 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)
7. Groenwold, R.H.: Informative missingness in electronic health record systems: the curse of knowing. *Diagnostic and prognostic research* **4**(1), 1–6 (2020)
8. Harutyunyan, H., Khachatryan, H., Kale, D.C., Ver Steeg, G., Galstyan, A.: Multi-task learning and benchmarking with clinical time series data. *Scientific data* **6**(1), 1–18 (2019)
9. Johnson, A.E., Pollard, T.J., Shen, L., Li-Wei, H.L., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Celi, L.A., Mark, R.G.: MIMIC-III, a freely accessible critical care database. *Scientific data* **3**(1), 1–9 (2016)
10. Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., Krishnan, D.: Supervised contrastive learning. *Advances in neural information processing systems* **33**, 18661–18673 (2020)
11. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014)
12. Le-Khac, P.H., Healy, G., Smeaton, A.F.: Contrastive representation learning: A framework and review. *Ieee Access* **8**, 193907–193934 (2020)
13. Lee, Y., Jun, E., Choi, J., Suk, H.I.: Multi-view integrative attention-based deep representation learning for irregular clinical time-series data. *IEEE Journal of Biomedical and Health Informatics* **26**(8), 4270–4280 (2022)
14. Li, J., Shang, J., McAuley, J.: Uctopic: Unsupervised contrastive learning for phrase representations and topic mining. *arXiv preprint arXiv:2202.13469* (2022)
15. Li, M., Li, C.G., Guo, J.: Cluster-guided asymmetric contrastive learning for unsupervised person re-identification. *IEEE Transactions on Image Processing* **31**, 3606–3617 (2022)
16. Li, R., Ma, F., Gao, J.: Integrating multimodal electronic health records for diagnosis prediction. In: *AMIA Annual Symposium Proceedings*. vol. 2021, p. 726. American Medical Informatics Association (2021)
17. Luo, Y., Cai, X., Zhang, Y., Xu, J., et al.: Multivariate time series imputation with generative adversarial networks. *Advances in neural information processing systems* **31** (2018)

18. Luo, Y., Zhang, Y., Cai, X., Yuan, X.: E2gan: End-to-end generative adversarial network for multivariate time series imputation. In: Proceedings of the 28th international joint conference on artificial intelligence. pp. 3094–3100. AAAI Press (2019)
19. Ma, L., Gao, J., Wang, Y., Zhang, C., Wang, J., Ruan, W., Tang, W., Gao, X., Ma, X.: Adacare: Explainable clinical health status representation learning via scale-adaptive feature extraction and recalibration. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 825–832 (2020)
20. Ma, L., Ma, X., Gao, J., Jiao, X., Yu, Z., Zhang, C., Ruan, W., Wang, Y., Tang, W., Wang, J.: Distilling knowledge from publicly available online emr data to emerging epidemic for prognosis. In: Proceedings of the Web Conference 2021. pp. 3558–3568 (2021)
21. Ma, L., Zhang, C., Wang, Y., Ruan, W., Wang, J., Tang, W., Ma, X., Gao, X., Gao, J.: Concare: Personalized clinical feature embedding via capturing the healthcare context. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 833–840 (2020)
22. McCombe, N., Liu, S., Ding, X., Prasad, G., Bucholc, M., Finn, D.P., Todd, S., McClean, P.L., Wong-Lin, K.: Practical strategies for extreme missing data imputation in dementia diagnosis. *IEEE journal of biomedical and health informatics* **26**(2), 818–827 (2021)
23. Mulyadi, A.W., Jun, E., Suk, H.I.: Uncertainty-aware variational-recurrent imputation network for clinical time series. *IEEE Transactions on Cybernetics* **52**(9), 9684–9694 (2021)
24. Ni, Q., Cao, X.: Mbgan: An improved generative adversarial network with multi-head self-attention and bidirectional rnn for time series imputation. *Engineering Applications of Artificial Intelligence* **115**, 105232 (2022)
25. Oh, E., Kim, T., Ji, Y., Khyalia, S.: Sting: Self-attention based time-series imputation networks using gan. In: 2021 IEEE International Conference on Data Mining (ICDM). pp. 1264–1269. IEEE (2021)
26. Pang, B., Li, Y., Zhang, Y., Peng, G., Tang, J., Zha, K., Li, J., Lu, C.: Unsupervised representation for semantic segmentation by implicit cycle-attention contrastive learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 2044–2052 (2022)
27. 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)
28. Pollard, T.J., Johnson, A.E., Raffa, J.D., Celi, L.A., Mark, R.G., Badawi, O.: The eicu collaborative research database, a freely available multi-center database for critical care research. *Scientific data* **5**(1), 1–13 (2018)
29. Shekhalishahi, S., Balaraman, V., Osmani, V.: Benchmarking machine learning models on multi-centre eicu critical care dataset. *Plos one* **15**(7), e0235424 (2020)
30. Shi, Z., Wang, S., Yue, L., Pang, L., Zuo, X., Zuo, W., Li, X.: Deep dynamic imputation of clinical time series for mortality prediction. *Information Sciences* **579**, 607–622 (2021)
31. Tan, Q., Ye, M., Yang, B., Liu, S., Ma, A.J., Yip, T.C.F., Wong, G.L.H., Yuen, P.: Data-gru: Dual-attention time-aware gated recurrent unit for irregular multivariate time series. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 930–937 (2020)
32. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. In: Advances in neural information processing systems. pp. 5998–6008 (2017)

33. Wang, W., Zhou, T., Yu, F., Dai, J., Konukoglu, E., Van Gool, L.: Exploring cross-image pixel contrast for semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7303–7313 (2021)
34. Wang, Y., Min, Y., Chen, X., Wu, J.: Multi-view graph contrastive representation learning for drug-drug interaction prediction. In: Proceedings of the Web Conference 2021. pp. 2921–2933 (2021)
35. Xu, D., Sheng, J.Q., Hu, P.J.H., Huang, T.S., Hsu, C.C.: 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* **25**(6), 2260–2272 (2020)
36. Yang, C., An, Z., Cai, L., Xu, Y.: Mutual contrastive learning for visual representation learning. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 3045–3053 (2022)
37. Yıldız, A.Y., Koç, E., Koç, A.: Multivariate time series imputation with transformers. *IEEE Signal Processing Letters* **29**, 2517–2521 (2022)
38. Yuan, X., Lin, Z., Kuen, J., Zhang, J., Wang, Y., Maire, M., Kale, A., Faieta, B.: Multimodal contrastive training for visual representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6995–7004 (2021)
39. Zang, C., Wang, F.: Scehr: Supervised contrastive learning for clinical risk prediction using electronic health records. In: Proceedings. IEEE International Conference on Data Mining. vol. 2021, pp. 857–866 (2021)
40. Zhang, Y., Zhou, B., Cai, X., Guo, W., Ding, X., Yuan, X.: Missing value imputation in multivariate time series with end-to-end generative adversarial networks. *Information Sciences* **551**, 67–82 (2021)