

A Multi-Task Imputation and Classification Neural Architecture for Early Prediction of Sepsis from Multivariate Clinical Time Series

Yale Chang*, Jonathan Rubin*, Gregory Boverman, Shruti Vij, Asif Rahman, Annamalai Natarajan, Saman Parvaneh

Philips Research North America, Cambridge, USA

* Authors contributed equally

Abstract

Early prediction of sepsis onset can notify clinicians to provide timely interventions to patients to improve their clinical outcomes. The key question motivating this work is: given a retrospective patient cohort consisting of multivariate clinical time series (e.g., vital signs and lab measurement) and patients' demographics, how to build a model to predict the onset of sepsis six hours earlier? To tackle this challenge, we first used a recurrent imputation for time series (RITS) approach to impute missing values in multivariate clinical time series. Second, we applied temporal convolutional networks (TCN) to the RITS-imputed data. Compared to other sequence prediction models, TCN can effectively control the size of sequence history. Third, when defining the loss function, we assigned custom time-dependent weights to different types of errors. Experiments on a real-world sepsis patient cohort from the PhysioNet Challenge 2019 demonstrate the effectiveness of the proposed approach.

1. Introduction

Sepsis is a life-threatening condition that puts 1.7 million lives at risk every year in the United States alone [1]. Early detection of sepsis in the intensive care unit would allow faster administration of antibiotic treatment and improve patient outcomes, as well as significantly reduce hospital expenses.

Automated systems that can accurately predict sepsis early based on available clinical data can allow faster treatment times and significant value for patients and care givers. However, there are various challenges associated with developing early prediction systems. First, predictions that are too early or false positive predictions can end up straining hospital resources, as well as being disadvantageous to patients. Second, predictions that are too late may reduce the impact of any administered intervention. Third, accuracy of predictions are affected by the nature

of clinical time series data, which is messy and consists of irregularly sampled and missing data points.

Previous works on sepsis prediction have employed time varying Cox proportional hazards models [2], gradient boosting [3] and Gaussian process RNNs [4]. To tackle the challenges associated with this problem, we build a multi-task neural architecture for early sepsis detection. First, we define a loss function to approximate the provided custom utility function [5], penalizing both too early and too late detections of sepsis. Second, we apply a recurrent imputation for time series (RITS) [6] model to impute missing values in multivariate clinical time series. Compared to other commonly-used imputation techniques, RITS can achieve significantly lower imputation error. By combining the RITS reconstruction loss with the utility loss, the imputed feature values are adaptive to the prediction task. Third, we train a TCN model using the RITS-imputed data to maximize the approximated utility objective. Experimental results on the PhysioNet Challenge 2019 patient cohort demonstrate the effectiveness of our proposed approach.

2. Imputation Network

The input of our neural architecture is a matrix, X , that consists of multivariate clinical time series (including vital signs, laboratory values and demographic information) that can contain missing values ($-$). Each row, in X records the value of a feature and the columns $x_1, x_2, \dots, x_t \in X$ record the value of each feature over time.

From X , we further derive two input matrices: 1) a matrix of masking vectors, m , that indicates whether the measurement is available (1) or missing (0); and 2) a δ matrix that records the time (in hours) since the last available feature measurement. Examples of input matrix X , masking vector m , and matrix δ are shown in Equations (1) – (4).

As the input vector, x_t , may contain missing values, it cannot be fed directly into a recurrent neural network architecture. Instead, we first process x_t using a RITS net-

work based on [6]. The RITS network uses a combination of *history-based estimation* and *feature-based estimation* to impute missing values. **History-based estimation** predicts missing feature values using the hidden state vector of a recurrent neural network. **Feature-based estimation** uses information about the available surrounding features at the current time, x_t , to make a prediction about missing variables.

$$X = \begin{matrix} & x_1 & x_2 & x_3 & \cdots & x_t \\ x^1 & \begin{bmatrix} 101 & 100 & - & 98 & 80 \end{bmatrix} \\ \vdots & \begin{bmatrix} 98.5 & - & - & 97 & 99 \end{bmatrix} \\ x^d & \begin{bmatrix} 37 & - & - & - & - \end{bmatrix} \end{matrix} \quad (1)$$

$$\mathbf{m}_t^d = \begin{cases} 0 & \text{if } x_t^d \text{ is unobserved} \\ 1 & \text{otherwise} \end{cases} \quad (2)$$

$$\delta_t^d = \begin{cases} 1 + \delta_{t-1}^d & \text{if } t > 1, \mathbf{m}_{t-1}^d = 0 \\ 1 & \text{if } t > 1, \mathbf{m}_{t-1}^d = 1 \\ 0 & \text{if } t = 1 \end{cases} \quad (3)$$

$$\mathbf{m} = \begin{bmatrix} 1 & 1 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \end{bmatrix} \delta = \begin{bmatrix} 0 & 1 & 1 & 2 & 1 \\ 0 & 1 & 2 & 3 & 1 \\ 0 & 1 & 2 & 3 & 4 \end{bmatrix} \quad (4)$$

A standard recurrent neural network formulation is shown in Equation (5). Here, $\sigma(\cdot)$ is some activation function, e.g. sigmoid, x_t, h_t are the input vector and recurrent hidden state vector, respectively, and W_h, U_h are learned network parameters.

$$\mathbf{h}_t = \sigma(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{U}_h \mathbf{x}_t + \mathbf{b}_h) \quad (5)$$

$$\hat{\mathbf{x}}_t = \mathbf{W}_x \mathbf{h}_{t-1} + \mathbf{b}_x \quad (6)$$

$$\mathbf{x}_t^c = \mathbf{m}_t \odot \mathbf{x}_t + (\mathbf{1} - \mathbf{m}_t) \odot \hat{\mathbf{x}}_t \quad (7)$$

$$\gamma_t = \exp\{-\max(0, \mathbf{W}_\gamma \delta_t + \mathbf{b}_\gamma)\} \quad (8)$$

$$\mathbf{h}_t = \sigma(\mathbf{W}_h [\mathbf{h}_{t-1} \odot \gamma_t] + \mathbf{U}_h [\mathbf{x}_t^c \odot \mathbf{m}_t] + \mathbf{b}_h) \quad (9)$$

$$\ell_t = \langle \mathbf{m}_t, \mathcal{L}_e(\mathbf{x}_t, \hat{\mathbf{x}}_t) \rangle \quad (10)$$

Equations (6) – (10) detail the operations required to perform imputation, based on a modified version of a recurrent neural architecture.

(6): Weight parameter matrices are learned and multiplied with the networks recurrent hidden state (at the previous time step) to estimate a new input vector with no missing information, \mathbf{x}_t .

(7): A complement vector, \mathbf{x}_t^c , is formed by combining both the available feature values with estimated information when the feature value is missing, where \odot represents element-wise multiplication.

(8): In addition to predicting missing information the neural architecture also applies a temporal decay, γ_t , to the

recurrent hidden state, based on the time gap since a feature value was measured, i.e. the δ matrix. The longer the time gap since the feature was measured results in more exponential decay being applied to the corresponding feature.

(9): The final hidden state vector, \mathbf{h}_t , is derived via a learned linear combination of *history-based estimation* (captured in the previously described steps) and *feature based estimation* that learns a further weight matrix, \mathbf{U}_h , to mix between available feature values at the present time step, where \circ is the concatenation operator. Equation (9) shows the final hidden state vector computed using a vanilla RNN formulation, whereas in practice an LSTM or GRU could be used to make an outcome prediction for the current time step.

(10): Finally, a loss function is applied (e.g. mean absolute error or mean squared error) between the estimated feature values and the true feature values (when available) that allows the network to learn its weight matrix parameters.

3. Sequence Prediction Network

Given input clinical time series ($\mathbf{x}_1^c, \dots, \mathbf{x}_t^c$), which are imputed through either RITS or forward-filling, a sequence prediction model $f(\cdot)$ maps the input sequence to an output sequence. In sepsis prediction, y_t represents the probability that the patient would have sepsis in six hours at the t -th hour.

$$\hat{y}_1, \dots, \hat{y}_t = f(\mathbf{x}_1^c, \dots, \mathbf{x}_t^c) \quad (11)$$

We select temporal convolutional network (TCN) [7] as the sequence prediction model due to its four desirable properties:

First, both the output sequence and each hidden layer have the same length with the input sequence, which is achieved by adding zero padding to subsequent layers. As a result, TCN can output the patient’s probability of having sepsis in six hours at any time during the ICU stay.

Second, TCN uses causal convolutions to avoid the leakage from future to the past.

Third, to enable TCN to utilize long history sequence, dilated convolutions are applied. By setting the dilation factor to be an exponential function of the network depth, the receptive fields of hidden layers grow exponentially w.r.t. the network depth. Therefore, the receptive fields can be flexibly adjusted by either changing 1) the number of hidden layers; or 2) the kernel size.

Fourth, the model can be trained in parallel and therefore faster than other sequence prediction models such as the recurrent neural network.

4. Utility Loss

We derive an approximation of the utility function defined in [5] and define it as the learning objective. For sepsis detection, the *optimal detection time* t^* is defined as six hours before the onset time. For a true sepsis patient, the prediction model would be rewarded 1 if it can correctly predict sepsis at the optimal detection time. Both too early prediction (earlier than six hours before the onset) or too late prediction (later than six hours before the onset or after the onset) would receive a reward less than 1. The reward (utility) would decrease as the prediction time moves further from the optimal detection time. The utility curves of true positive U_{TP} and false negative U_{FN} are illustrated in the top subfigure of Figure 1 (from [5]).

For a non-sepsis patient, there would be slight penalty for falsely predicting sepsis. There would be no reward or penalty for correctly identifying the patient as non-sepsis. The utility curves of false positive U_{FP} and true negative U_{TN} are shown in the bottom subfigure in Figure 1 (from [5]).

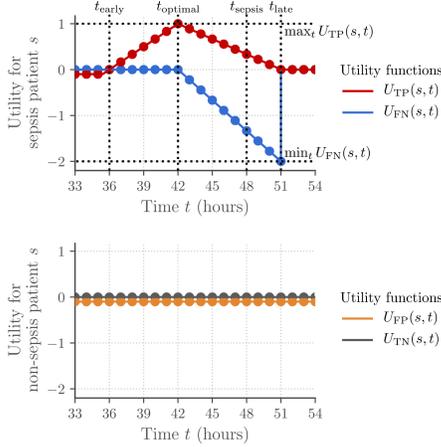


Figure 1. The top subfigure shows the utility curve of true positive U_{TP} and false negative U_{FN} . The bottom subfigure shows the utility curve of false positive U_{FP} and true negative U_{TN} .

Specifically, the utility curve can be expressed as follows

$$U_{TP}(t) = -0.05\mathbb{I}[t < t^* - 6] + (t - t^* + 6)/6 \quad (12)$$

$$\mathbb{I}[t^* - 6 \leq t \leq t^*] + (9 - t + t^*)/9\mathbb{I}[t^* < t \leq t^* + 9]$$

$$U_{FN}(t) = 0\mathbb{I}[t < t^*] - 2(t - t^*)/9\mathbb{I}[t^* \leq t \leq t^* + 9] \quad (13)$$

$$U_{FP}(t) = -0.05; \quad U_{TN}(t) = 0 \quad (14)$$

For each of N patients, given the sepsis label s_n , optimal detection time t_n^* , and predicted probability of sepsis in six hours at the t -th hour during the patient's ICU stay

p_{nt} , we define the learning objective as maximizing the approximated utility function

$$\theta^* = \arg \max_{\theta} \sum_{n=1}^N \sum_{t=1}^{T_n} \mathbb{I}[s_n = 1] (p_{nt} U_{TP}(t) + (1 - p_{nt}) U_{FN}(t)) + \mathbb{I}[s_n = 0] (p_{nt} U_{FP}(t) + (1 - p_{nt}) U_{TN}(t)) \quad (15)$$

where p_{nt} is modeled as TCN, θ consists of parameters of TCN.

5. Results

The dataset consists of multivariate clinical time series (vital signs, lab measurements) and static variables (demographics, unit admission information) from 40338 patient ICU stays. The details of the cohort are described in [5].

To train the imputation network, we randomly split the dataset into 10 folds of equal size and apply 10-fold *nested* cross validation. Among these 10 folds, 8 folds are used as training set, 1 fold is used as validation set and the remaining 1 fold is used as the test set. When the fold i ($i = 1, \dots, 9$) is used as the test set, fold $i - 1$ is used as the validation set. For test fold 0, fold 9 is used as the validation set.

When training the imputation network, we regularize the reconstruction loss with the utility loss to make the imputed values adaptive to the prediction. Both the network architecture and hyperparameters are tuned by optimizing the loss function on the validation set. The resulting LSTM cell contains one hidden layer of size 64 with dropout rate 0.5.

After imputing missing values, we train a TCN to maximize the approximated utility function. The input features consist of 1) imputed feature values output by the imputation network; and 2) missingness indicator variables of the raw vital signs and lab measurements. We use the same train-validation-test split as training the imputation network. After tuning the TCN architecture and hyperparameters, the optimal TCN consists of one hidden layer with 100 kernels, where the size of each kernel is 5.

We compared the following approaches in terms of test normalized utility:

- **RITS**: apply the trained RITS model to test set and compute utility.
- **FF-TCN**: train TCN on forward-filled time series and missingness indicator variables
- **RITS-TCN**: train TCN on RITS-imputed time series and missingness indicator variables

For these three approaches, we show their test normalized utility scores of 10-fold cross validation in table 1.

First, comparing RITS-TCN and RITS, RITS-TCN consistently outperforms RITS on each test fold in achieving

higher test utility score. Therefore, it’s necessary to train a separate prediction network even though RITS can combine the reconstruction loss (for imputation) and the utility loss (for prediction) in its objective.

Second, comparing RITS-TCN and FF-TCN, RITS-TCN outperforms FF-TCN on 8 of 10 test folds, indicating the advantage of applying RITS to impute missing values instead of simple forward filling.

Third, comparing different folds within each model, the test utility scores have large variance, indicating the necessity of building ensemble models. This is because the ensemble model can reduce variance, leading to an improved utility score on the test set.

Table 1. 10-fold cross validation normalized utility scores of RITS, FF-TCN and RITS-TCN

Fold	RITS	FF-TCN	RITS-TCN
0	0.373	0.396	0.397
1	0.315	0.357	0.354
2	0.395	0.393	0.415
3	0.352	0.414	0.425
4	0.366	0.356	0.382
5	0.361	0.398	0.392
6	0.336	0.368	0.376
7	0.387	0.388	0.416
8	0.371	0.374	0.387
9	0.401	0.426	0.430
Mean (Std)	0.366 (0.025)	0.387 (0.022)	0.397 (0.024)

We submit 3-RITS-TCN-Ensemble, an ensemble of 3 RITS-TCN models trained using fold 3, 7, 9 as the test set respectively. The ensemble model would predict sepsis if at least one model would predict sepsis. Evaluated on a subset of the hidden challenge test set, our model achieved a normalized utility score of 0.411. We also compare the performance of 3-RITS-TCN-Ensemble against RITS, FF-TCN, RITS-TCN on the subset of the hidden test set. The results are shown in Table 2.

Table 2. Test utility scores of RITS, FF-TCN, RITS-TCN, 3-RITS-TCN-Ensemble on subset of the hidden test set

Model	Normalized Utility
RITS	0.356
FF-TCN	0.389
RITS-TCN	0.399
3-RITS-TCN-Ensemble	0.411

Similar to the evaluation on the in-house test sets, RITS-TCN outperforms both FF-TCN and RITS. Furthermore, 3-RITS-TCN-Ensemble outperforms RITS-TCN, confirming our hypothesis that the ensemble model can lead to improved test utility due to the reduction of predictive variance compared to a single RITS-TCN.

6. Conclusions and Future Work

In this work, we build a sepsis prediction model by applying RITS to impute missing values in multivariate clinical time series followed by TCN to maximize the approximated utility objective. Experimental results on the PhysioNet Challenge 2019 sepsis cohort demonstrate the effectiveness of our proposed approach.

Both RITS and TCN are black-box models. In order to incorporate them into clinicians’ workflows to provide clinical decision supports, we still need to provide interpretations for model predictions. In the future work, we would work on interpreting the model through the attention mechanism.

The performance of the ensemble model can be further improved by increasing the number of RITS-TCN models. Other prediction models, such as XGBoost can also be added to the ensemble to increase the diversity of base models, leading to higher test utility scores.

References

- [1] <https://www.cdc.gov/sepsis/datareports/index.html>.
- [2] Nemati S, Holder A, Razmi F, Stanley MD, Clifford GD, Buchman TG. An interpretable machine learning model for accurate prediction of sepsis in the icu. *Critical care medicine* 2018;46(4):547–553.
- [3] Barton C, Chettipally U, Zhou Y, Jiang Z, Lynn-Palevsky A, Le S, Calvert J, Das R. Evaluation of a machine learning algorithm for up to 48-hour advance prediction of sepsis using six vital signs. *Computers in biology and medicine* 2019; 109:79–84.
- [4] Futoma J, Hariharan S, Sendak M, Brajer N, Clement M, Bedoya A, O’Brien C, Heller K. An improved multi-output gaussian process rnn with real-time validation for early sepsis detection. *2nd Machine Learning for Healthcare Conference MLHC 2017*.
- [5] Reyna MA, Josef C, Jeter R, Shashikumar SP, M. Brandon Westover MB, Nemati S, Clifford GD, Sharma A. Early prediction of sepsis from clinical data: the PhysioNet/Computing in Cardiology Challenge 2019. *Critical Care Medicine* 2019;In press.
- [6] Cao W, Wang D, Li J, Zhou H, Li L, Li Y. Brits: Bidirectional recurrent imputation for time series. In *Advances in Neural Information Processing Systems*. 2018; 6775–6785.
- [7] Bai S, Kolter JZ, Koltun V. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv180301271* 2018;.

Address for correspondence:

Yale Chang / Jonathan Rubin
 2 Canal Park, 3rd floor, Cambridge, MA 02141, USA
yale.chang@philips.com / jonathan.rubin@philips.com