

Predicting COVID-19 Disease Progression and Patient Outcomes based on Temporal Deep Learning

Chenxi Sun

Peking University

Shenda Hong

Peking University

Moxian Song

Peking University

Hongyan Li

Peking University

Zhenjie Wang (✉ zhenjie.wang@pku.edu.cn)

Peking University <https://orcid.org/0000-0003-1717-5790>

Technical advance

Keywords: COVID-19, Disease progression, Outcome early prediction, Irregularly sampled time series, Time-aware long short-term memory

Posted Date: November 6th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-44308/v2>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on February 8th, 2021. See the published version at <https://doi.org/10.1186/s12911-020-01359-9>.

1 **Predicting COVID-19 Disease Progression and Patient Outcomes based on**
2 **Temporal Deep Learning**

3

4 Chenxi Sun^{1,2}, Shenda Hong^{3,4}, Moxian Song^{1,2}, Hongyan Li^{1,2}, and Zhenjie Wang^{5,*}

5

6 ¹School of Electronics Engineering and Computer Science, Peking University, Beijing,
7 People's Republic of China.

8 ²Key Laboratory of Machine Perception (Ministry of Education), Peking University, Beijing,
9 People's Republic of China.

10 ³National Institute of Health Data Science, Peking University, Beijing, People's Republic of
11 China.

12 ⁴Institute of Medical Technology, Health Science Center of Peking University, Beijing, People's
13 Republic of China.

14 ⁵Institute of Population Research, Peking University, Beijing, People's Republic of China.

15

16 *Corresponding Author: Zhenjie Wang, Peking University, No.5 Yiheyuan Road, Beijing
17 100871, People's Republic of China. Email: zhenjie.wang@pku.edu.cn

18

1 **Abstract**

2 **Background:** The coronavirus disease 2019 (COVID-19) pandemic has caused health concerns
3 worldwide since December 2019. From the beginning of infection, patients will progress
4 through different symptom stages, such as fever, dyspnea or even death. Identifying disease
5 progression and predicting patient outcome at an early stage helps target treatment and resource
6 allocation. However, there is no clear COVID-19 stage definition, and few studies have
7 addressed characterizing COVID-19 progression, making the need for this study evident.

8 **Methods:** We proposed a temporal deep learning method, based on a time-aware long short-
9 term memory (T-LSTM) neural network and used an online open dataset, including blood
10 samples of 485 patients from Wuhan, China, to train the model. Our method can grasp the
11 dynamic relations in irregularly sampled time series, which is ignored by existing works.
12 Specifically, our method predicted the outcome of COVID-19 patients by considering both the
13 biomarkers and the irregular time intervals. Then, we used the patient representations, extracted
14 from T-LSTM units, to subtype the patient stages and describe the disease progression of
15 COVID-19.

16 **Results:** Using our method, the accuracy of the outcome of prediction results was more than
17 90% at 12 days and 98%, 95% and 93% at 3, 6, and 9 days, respectively. Most importantly, we
18 found 4 stages of COVID-19 progression with different patient statuses and mortality risks. We
19 ranked 40 biomarkers related to disease and gave the reference values of them for each stage.
20 Top 5 is Lymph, LDH, hs-CRP, Indirect Bilirubin, Creatinine. Besides, we have found 3
21 complications - myocardial injury, liver function injury and renal function injury. Predicting

1 which of the 4 stages the patient is currently in can help doctors better assess and cure the
2 patient.

3 **Conclusions:** To combat the COVID-19 epidemic, this paper aims to help clinicians better
4 assess and treat infected patients, provide relevant researchers with potential disease
5 progression patterns, and enable more effective use of medical resources. Our method predicted
6 patient outcomes with high accuracy and identified a four-stage disease progression. We hope
7 that the obtained results and patterns will aid in fighting the disease.

8

9 **Keywords:** COVID-19, Disease progression, Outcome early prediction, Irregularly sampled
10 time series, Time-aware long short-term memory.

11

1 **Background**

2 Coronavirus disease 2019 (COVID-19) outbreaks have caused health concerns worldwide since
3 December 2019; the disease was declared a pandemic by the World Health Organization (WHO)
4 on 11 March 2020 [1]. Over seven million cases of COVID-19 have been reported worldwide,
5 including more than 400,000 deaths (as of 15 June 2020) [2]. Even though the disease has been
6 controlled in certain countries, the WHO director warns the pandemic is still ‘Speeding Up’
7 [25]. Because of its sudden onset, many hospitals are still facing medical resource shortages.
8 For example, news in [26] reported a lack of medical resources in New Delhi. In [27], Arizona
9 has experienced record-high hospital capacity as coronavirus cases climb. A reasonable
10 allocation of resources according to patient condition is needed.

11 The solution to this problem involves determining the stages of disease progression by
12 subtyping and predicting the outcome of COVID-19 patients. Then, targeted treatment and
13 medical resource allocation can be carried out for patients in different stages. Recent studies [3-
14 5,14-16] have used statistical methods to analyze COVID-19 progress by inpatient symptoms.
15 However, different statistical results were obtained by considering different patient groups and
16 different symptoms. At present, there is no clear division of the stages of COVID-19
17 progression.

18 Longitudinal disease analysis is the key to understanding disease progression, designing
19 prognoses and developing early diagnostic tools. The time dynamics of disease can provide
20 more information than static symptom observation [46]. Considering the complex patient states,
21 the amount of interventions and the real-time requirement, the data-driven machine learning
22 approaches by learning from electronic health records are the desiderata to help clinicians [44].

1 Many existing works have used machine learning methods for COVID-19 prediction tasks.
2 We have summarized them in Table 1. For example, in most method of [45] and in [1,32,36-
3 38,40,42], authors used non-deep learning methods, such as k-NN, LR, Cox, SVM and DT to
4 classify CT/X-ray images and predict the outcomes of COVID-19 patients. However, in terms
5 of prediction accuracy, non-deep learning is not as good as deep learning methods. Deep
6 learning methods can train the parameters with complex nonlinearity to learn the data structures
7 and have achieved state-of-the-art in many medical prediction tasks [19,20,47]. Thus, many
8 current works apply deep learning methods for COVID-19 prediction tasks [28-36]. However,
9 these methods either use the simple multi-layer perceptron for predicting or use the
10 convolutional structures for image classification. Both the above methods ignored the temporal
11 development of patient's status. In the real-world patient records, except for the basic
12 information, vital signs, test values and diagnoses are both time series, especially for the blood
13 samples of COVID-19 patients, the data we used in this paper.

14 Recently, a deep learning method, recurrent neural network (RNN) [18] can efficiently
15 model temporal sequences. It uses recursion in the direction of sequence evolution to learning
16 the relations among past, present and future. But the basic RNN has the long-term dependency
17 problems [7]. Meanwhile, RNN only process uniformly distributed longitudinal data while
18 COVID-19 patient blood samples are distributed nonuniformly with irregular time intervals
19 between observations. Thus, a method that can model this irregular time series of COVID-19
20 patients is needed.

21 In this paper, we retrospectively analyzed the blood samples of 485 patients from the region
22 of Wuhan, China. The medical records collected with standard case report forms, including

1 epidemiological, demographic, clinical, laboratory and mortality outcome information, from an
2 online open dataset under an MIT license. We applied a temporal deep learning method Time-
3 aware Long Short-term Unit (T-LSTM) to model the irregular time series of COVID-19 patients.
4 T-LSTM can predict the mortality with more than 98% accuracy before 3 days. Meanwhile, we
5 have discovered four stages of COVID-19 patients. According to the different stages, we gave
6 the analysis of the patient's state and found the related biomarkers and complications.

7

8 **Methods**

9 In this section, we first introduce the COVID-19 dataset and the data preprocessing process.
10 Then, we describe the methods for mortality prediction and disease progression in detail.

11

12 **Dataset description**

13 Blood index values can reflect a COVID-19 patient's physical condition [15]. COVID-19
14 patients' blood samples were collected between 10 January and 18 February 2020 at Tongji
15 Hospital of Tongji Medical College, Huazhong University of Science and Technology, Wuhan,
16 China [6]. The dataset contains 80 characteristics from 375 patients with 6120 records as a
17 training set and 110 patients with 757 records as a test set. A case of sample is shown in Figure
18 1. It draws lines of the time series of LHD, lymph and hs-CRP of a 70-year-old female patient
19 during hospitalization. We can see the time intervals between two observations are irregular,
20 which could be a few minutes or even days.

21 The detailed statistical information of demographic and 74 clinical laboratory test features
22 is listed Table 2. For example, in the dataset, the average age of patients is 58.83, the survival

1 rate is 53.6% and the ratio of male to female is about 1.5:1. We also list the range and mean
 2 value of each feature. In Figure 1, we display the distributions of some features (age, gender,
 3 LHD, lymph and hs-CRP) of survival class (0) and death class (1).

4 This COVID-19 blood test data is publicly available at
 5 https://github.com/HAIRLAB/Pre_Surv_COVID_19.

6

7 **Dataset preprocessing**

8 First, we attempted to find a suitable time measurement granularity. In the raw dataset, the
 9 lengths of sequences are unequal and different sampling times result in missing data, with an
 10 85% missing rate (mr in Equation 1) on average. The presence of vacancies has a large impact
 11 on data quality, resulting in unstable predictions and other unpredictable effects [17]. We used
 12 3 days as the basic sampling interval, reducing the average mr below 30%. The time series
 13 length of raw data, the average missing rate and the missing rate for each feature are shown in
 14 Figure 1.

$$15 \quad mr = \frac{\# \text{ of records}}{\text{sequence length}} \quad (1)$$

16 Meanwhile, for feature selection, using all 74 laboratory test features is unrealistic. To
 17 address the high missing rate, repeated features and collection difficulties, we considered three
 18 key features: lactic dehydrogenase (LDH), lymphocytes (lymph) and high-sensitivity C-
 19 reactive protein (hs-CRP). These features contain specific research biomarkers of COVID-19
 20 patients [6] and can be easily collected in any hospital. Considering that only three features may
 21 not achieve high prediction accuracy, we also select 40 features (listed in Table 7) with missing
 22 rate less than 30% for comparative experiment.

1

2 **T-LSTM**

3 Recurrent neural networks (RNNs) [18] (the first structure in Figure 2) are deep network
 4 architectures designed to model temporal sequences. They take sequence data as input,
 5 recursion occurs in the direction of sequence evolution, and all units are chained together. In
 6 basic RNN (the second structure in Figure 2), the current state h_t is affected by the previous
 7 state h_{t-1} and the current input x_t and is described as $h_t = \sigma(Wx_t + Uh_{t-1} + b)$, where σ
 8 is an activation function, and W , U and b are learnable parameters. Long Short-Term
 9 Memory (LSTM) [7] (the third structure in Figure 2) is a variant of RNN that is adept at solving
 10 long-term dependency problems. A standard LSTM unit consists of a forget gate, an input gate,
 11 a memory cell and an output gate.

12 However, RNNs only process uniformly distributed longitudinal data by assuming that the
 13 sequences have an equal distribution of time differences. COVID-19 patient blood samples are
 14 distributed nonuniformly. For example, the time gap between two sequential records could be
 15 hours or days. Time-aware Long Short-Term Memory (T-LSTM) [8] (the fourth structure in
 16 Figure 2) incorporates the elapsed time information into LSTM. It applies a memory discount
 17 to capture the irregular temporal dynamics. T-LSTM can be formulated as:

18 (2)

$$19 \quad C_{t-1}^S = \tanh(W_d C_{t-1} + b_d) \quad \text{Short-term memory}$$

$$20 \quad \hat{C}_{t-1}^S = C_{t-1}^S * g(\Delta_t) \quad \text{Discounted short-term memory}$$

$$21 \quad C_{t-1}^T = C_{t-1} - C_{t-1}^S \quad \text{Long-term memory}$$

1	$C_{t-1}^* = C_{t-1}^T - \hat{C}_{t-1}^S$	Adjusted previous memory
2	$f_t = \sigma(W_f x_t + U_f h_{t-1} + b_f)$	Forget gate
3	$i_t = \sigma(W_i x_t + U_i h_{t-1} + b_i)$	Input gate
4	$\tilde{C}_t = \tanh(W_c x_t + U_c h_{t-1} + b_o)$	Candidate memory
5	$C_t = f_t * C_{t-1}^* + i_t * \tilde{C}_t$	Current memory
6	$o_t = \sigma(W_o x_t + U_o h_{t-1} + b_o)$	Output gate
7	$h_t = o_t * \tanh(C_t)$	Current hidden state

8 Here, we use a log calculation for the elapsed time function. Δ_t describes the time gap
9 between two records at two sequential time steps t and $t - 1$. T_t is the actual time at time
10 step t .

$$11 \quad g(\Delta_t) = \frac{1}{\log(e + \Delta_t)} \quad (3)$$

$$12 \quad \Delta_t = T_t - T_{t-1}$$

13

14 **Analysis strategy**

15 We first describe the two tasks in this study and then introduce the specific methods. The whole
16 method process is shown in Figure 3.

17 **Task 1 (Outcome prediction)** A set of labeled patient data is represented as $\mathcal{D} = \{(x_i, c_i) \in$
18 $(X, C) | i = 1, \dots, n\}$. X is a time series set of patients, where $x_i = \{x_i^t | t = 1, \dots, t_{onset}\}$
19 represents a patient's records over t time steps; specifically, x_i^t is multivariate data, and each
20 dimension is a clinical record represented by a numeric vector. $C \in \{0, 1\}$ is the outcome,
21 where class 0 means death and class 1 means survival. The outcome prediction task aims to

1 *predict patient outcomes by the prediction function $f: X \rightarrow C$.*

2 **Task 2 (Temporal patient subtyping / Disease progression mining)** *The goal is to find patient*
 3 *groups $G = \{g_i | i = 0, \dots, m\}$ with similar feature representation $R = \{r_i^t | i = 0, \dots, n; t =$
 4 $0, \dots, t_{onset}\}$. r_i^t is the representation of clinical record x_i^t at time t . Then, the patient groups*
 5 *G distributed over time are used to analyze the stages of disease progression.*

6 In COVID-19 patient outcome prediction task, T-LSTM is used to handle patient record
 7 sequences and then make the prediction. The process is displayed in the proposed method of
 8 Figure 2, in the lower gray area.

9 For a patient i , the input of T-LSTM at time step t is a three-dimensional feature vector
 10 $x_i^t = [v_{LDH}, v_{lymphocytes}, v_{hs-CRP}]$ with time gap Δ_t . The output is the state representation
 11 s_i at the last time step. We apply this outcome prediction task as a binary classification task,
 12 with two classes: death and survival.

13 The cross-entropy [43] is mainly used to measure the difference between two probability
 14 distributions. We expect our predicted distribution of patient outcomes to be closer to the true
 15 distribution. Thus, we use the cross-entropy loss function in Equation 4. Besides, when using
 16 sigmoid active function, this loss can avoid the reduced learning rate causing by traditional
 17 mean square error loss when gradient decreases.

$$18 \quad L = L_{CE}(C, \hat{C}) = -\sum_x p(x) \log q(x) = -\sum_{i=1}^n \hat{c}_i \log c_i + (1 - \hat{c}_i) \log(1 - c_i) \quad (4)$$

19 $p(x)$ is the prior probability (true label vector) and $q(x)$ is the prediction probability
 20 (predicted results vector). Correspondingly, \hat{C} is the real class of input data, and C represents
 21 the prediction class.

22 In COVID-19 patient disease progression task, temporal patient subtyping can uncover the

1 dynamic characteristics of diseases by significantly nuanced subtyping, which leads to the
 2 potential stages of disease progression. We addressed the issue by building a time stage
 3 reference and providing a low-dimensional representation of each subject, encoding his or her
 4 position with respect to this reference.

5 The method structure is displayed in the upper gray area of proposed method in Figure 2.
 6 It has 4 steps: 1) Acquisition of patient representation r^t . We used the hidden state h_t ,
 7 extracted from every T-LSTM unit, as the patient's representation r^t at time step t . 2)
 8 Dimension reduction of r^t . For better demonstration, we used the t-distributed Stochastic
 9 Neighbor Embedding (t-SNE) [9] method to reduce these high-dimensional vectors r^t into
 10 two dimensions. 3) Obtaining the patient group set G . As prior information about the patient
 11 groups was not available, we acquired patient groups by applying an unsupervised clustering
 12 method, the Density-Based Spatial Clustering of Applications with Noise (DBSCAN) [10], on
 13 r^t . 4) Analysis of G and stages of disease progression. The mortality rate, MR , and the
 14 average time distance, TD , were calculated as the analysis criteria:

$$15 \quad MR = \frac{\# \text{ of death}}{\# \text{ of patients}} \quad (5)$$

$$16 \quad TD = \frac{1}{|g_k|} \sum_{x_i^t \in g_k} (T_{t_{onset}} - T_t) \quad (6)$$

17

18 **Evaluation metrics**

19 The prediction results were evaluated by assessing the area under the curve of the Receiver
 20 Operating Characteristic (AUC-ROC). The ROC is a curve of the True Positive Rate (TPR) and
 21 the False Positive Rate (FPR). TN, TP, FP and FN represent true positives, true negatives, false
 22 positives and false negatives, respectively.

$$1 \quad TPR = \frac{TP}{TP+FN} \quad (7)$$

$$2 \quad FPR = \frac{FP}{TN+FP} \quad (8)$$

3 The patient groups obtained by unsupervised clustering were evaluated by the Calinski-
 4 Harabaz Index (CH), which measures the covariance of data within a class and between classes.
 5 A larger CH value indicates a better clustering performance. In Equation 9, m is the number
 6 of data and k is the number of groups. B_k and W_k respectively represent the covariance
 7 matrices between groups and within groups.

$$8 \quad CH = \frac{\text{tr}(B_k) \frac{m-k}{k-1}}{\text{tr}(W_k)} \quad (9)$$

9 When we get the stages of COVID-19 patients, we used Kullback-Leibler Divergence (KL
 10 divergence) to analyze patient characteristics through each laboratory test feature. KL
 11 divergence can measure the degree of difference between two probability distributions. For
 12 each feature, we first establish the Gaussian distribution at each stage. Then, we calculate the
 13 average KL divergence of the distribution of adjacent stages. If the average KL divergence of a
 14 feature is large, it more likely is a biomarker to distinguish different stages. The basic KL
 15 divergence of distribution $p(X)$ and $q(X)$ and the KL divergence of two univariate Gaussian
 16 distributions are in Equation 10 and 11.

$$17 \quad KL(p(X)||q(X)) = \sum_{x_i \in X} p(x_i) \log \frac{p(x_i)}{q(x_i)} \quad (10)$$

$$18 \quad KL(\mathcal{N}(\mu_1, \sigma_1^2)||\mathcal{N}(\mu_2, \sigma_2^2)) = \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2} \quad (11)$$

19 For measure and evaluate each feature, we use the average KL divergence (Average KL)
 20 between neighbor stages g_i, g_{i+1} .

$$21 \quad \text{Average KL} = \frac{1}{m} \sum_{i=0}^{m-1} KL_{g_i, g_{i+1}} \quad (12)$$

22

1 **Results**

2 We used the records of 375 patients as a training set; the ratio of the training set to the
3 verification set was 0.8:0.2. The records of 110 patients made up the test set. This experiment
4 was conducted on 5-fold cross-validation. The code implementation is publicly available at
5 <https://github.com/scxhhh/COVID-19>.

6

7 **Baselines**

8 We use the related works summarized in Table 1 as comparison methods. Related works are
9 divided into non-deep learning methods and deep learning methods. We use Cox [32], k-NN
10 [42], SVM [36], DT [1], BPNN [28], PNN [29], RNN, LSTM and T-LSTM for COVID-19
11 mortality prediction. T-LSTM is our method.

12

13 **Outcome prediction results.**

14 Table 3 shows the results of COVID-19 mortality prediction using baselines. The AUC-ROC is
15 evaluated at 0, 3, 6, 9, 12, 15, and 18 days early. Here, the results are obtained when the patient's
16 representations are 64 dimensional. The results indicate that our method T-LSTM performed
17 better than all of baselines no matter how early before the onset times of patients. More precisely,
18 using T-LSTM, the outcome prediction accuracy is above 90% at 12 days early and is
19 approximately 97% accurate when predicting 3 days before the disease outcome. More detailed
20 results of train, validation and test sets using T-LSTM are listed in Table 4.

21 The first four figures in Figure 3 are the visualizes of prediction results. The first two
22 figures are the AUC-ROC of prediction results of baselines and T-LSTM in different earliness.

1 The third figure is the changes of prediction accuracy and cross-entropy loss when training the
2 model. The fourth figure represents the relation of patient representation dimension and AUC-
3 ROC of prediction using T-LSTM. Too few dimensions lead to incomplete feature learning,
4 while too many dimensions lead to redundant calculations and easy over-fitting. Considering
5 result accuracy, computational complexity and ease of representation use in the following task,
6 we decided to use 64 dimensional vectors to represent patients.

7 Further, we also select 40 features (listed in Table 7) as the input of T-LSTM for
8 comparative experiment. The results in Table 5 indicate that learning a large number of patient
9 characteristics does not necessarily contribute to COVID-19 patient mortality prediction task.
10 The three biomarkers, LDH, lymph and hs-CRP can make the results better. The AUC-ROC of
11 using 3 features is 3% higher than using 40 features on average.

12

13 **Disease progression results.**

14 By implementing the four steps of disease progression mining, we obtained the 4 stages in both
15 the death class (critical) and the survival class (general), shown in Figure 4.

16 For better visualization, we reduced the dimension of the patient's representation vector,
17 the fifth figure in Figure 3 is the dimension reduction effect. We chose 2 dimensions due to low
18 representation loss and clear observation. Besides, the DBSCAN clustering effect evaluated by
19 the CH index is shown in the sixth and seventh figures in Figure 3. Different clustering effects
20 can be obtained by changing the cluster radius parameter ϵ . The best CH index values for the
21 death class and the survival class are 680.07 and 44.24, respectively.

22 In this case, both classes have four groups. Four stages of COVID-19 patients are shown

1 in Figure 4. For each stage, we calculate the mortality rate MR and the average time distance
2 TD . For the death class, MR increases over stages and is 100% at stage 4. For the survival
3 class, MR decreases over stages and is 0% in stage 4. TD in both classes decreases, meaning
4 that the 4 stages are distributed over time. Meanwhile, as the CH index of the survival class is
5 higher than that of the death class, the survival class stages are relatively loosely distributed.

6 In Figure 4, the first clustering is obtained by using biomarkers directly and shows that
7 reasonable stages could not be found. In the first clustering, no stage is clustered in the death
8 class and the 2 stages in the survival class have similar mortality rates and no time difference,
9 as the shade of blue indicates. However, using our method, different stages have obvious
10 differences, such as the data point color deepening with the stages. Meanwhile, as shown in the
11 two insets, the class boundary is clearer based on our method.

12 Meanwhile, we calculated the mean values of 40 laboratory test features in each stage, the
13 feature values vary with stages. Table 6 lists 10 of these features - Lymph, LDH, hs-CRP,
14 Indirect Bilirubin, Creatinine, INR, Serum Sodium, eGFR, Serum Chlorine and Albumin. The
15 changes of values through 4 stages are visualized in Figure 5. Under different classes, the trends
16 of features are different.

17 Further, we calculated the average KL divergence between adjoint stages of each features
18 in 40 clinical laboratory tests data. We ranked the average KL values. The higher the ranking,
19 the better the biomarkers can be used to distinguish different stages. The top five are: Lymph,
20 LDH, hs-CRP, Indirect Bilirubin and Creatinine. Among them, Lymph is the best biomarker to
21 distinguish stages.

22

1 **Discussion**

2 In recent years, deep learning (DL) technology has been widely used because of its superior
3 performance in various medical applications [19, 20], such as medical image recognition [21]
4 and medication recommendations [22]. In the past year, the spread of COVID-19 has had a
5 peripheral effect on the global economy and health. Therefore, we expect to combine DL
6 methods to study and fight COVID-19.

7 The states of COVID-19 patients in hospital are dynamic time sequence processes. In
8 addition to the basic information of patients, the vital signs, diagnoses and other lab tests are all
9 time series. Existing many works [28-42,45] have achieved good results for COVID-19
10 prediction tasks. But they paid little attention to analyze and model the characteristics of
11 COVID-19 patients' time series. Dynamic time series modeling can grasp the relationship
12 between historical observations and current observations, and learn the potential development
13 mode of sequence, which is conducive to more accurate prediction and representation. Besides,
14 we have found that the time series of COVID-19 patients is irregularly sampled - Different time
15 intervals exist in adjacent observations. Every possible test is not regularly measured during an
16 admission. When a certain symptom worsens, corresponding variables are examined more
17 frequently; when the symptom disappears, the corresponding variables are no longer examined.
18 These time intervals will add a time sparsity factor when the intervals between observations are
19 large [44]. Therefore, it is necessary not only to deal with time series, but also to deal with
20 irregular time series according to the characteristics of COVID-19 patients.

21 This study has two basic contributions. First, we can predict patient outcomes with higher
22 accuracy than baselines. The method can effectively predict whether the infected patient will

1 die or survive 12 days prior to disease outcome with over 90% accuracy. The prediction
2 accuracies at 3, 6, and 9 days prior are 98%, 95% and 93%, respectively. Second, we identified
3 four stages of COVID-19 progression. The stages are closely related to mortality and time of
4 illness and can help analyze the status of infected patients. Through the analysis of each stage,
5 we ranked 40 biomarkers according to the degree of correlation with COVID-19.

6 Currently, certain studies have identified suitable predictive biomarkers, such as the 3
7 biomarkers in [6], which are regarded to have a significant impact on patient mortality. If a
8 doctor predicts survival or death only by observing the biomarkers and using a threshold, the
9 accuracy is at or below 80% for early predictions. However, the clinical reference value of
10 inaccurate results is very low [11, 12]. Finding the appropriate prediction-related biomarkers is
11 important, but it is equally important to create an appropriate and high-accuracy prediction
12 method. The DL method has better performance and the time-aware aspect enables higher
13 accuracy, as shown in Table 3. Decision trees [6] are easy to explain, but do not have good
14 accuracy. T-LSTM can better grasp the relationship between time than basic LSTM. It not only
15 improves the prediction accuracy but also contributes to the identification of disease stages, as
16 the disease progression is distributed over time.

17 However, there are some concerns about the use of DL methods in the high-risk tasks of
18 healthcare. First, we recognize that it may be risky to apply predictive methods directly to
19 clinical practice [13]. DL methods may be assistive tools for doctors but not used to make
20 decisions directly. It is challenging for doctors to make optimal decisions; a data-driven, high-
21 accuracy prediction method could help. Second, the DL method is troubled by poor
22 interpretability [23], and clinical settings prefer interpretable models. Our approach has this

1 flaw, as do other black-box models. We feel it is best to keep a complex black-box model to
2 gain 3 percent higher accuracy, compared with interpretable methods [6]. We also wish to
3 provide interpretable results to help doctors better understand the model's result. Thus, we
4 ranked 40 biomarkers according to the degree of correlation with COVID-19 in disease
5 progression task.

6 To help relevant researchers better study COVID-19 patients, we have uncovered the
7 disease progression of COVID-19 patients and obtained 4 stages. This interesting finding
8 cannot be distinguished simply by the value of biomarkers, as shown as the comparison of two
9 clustering results in Figure 4. Our method found stages closely related to mortality and time of
10 illness. This shows that the DL method can explore new patterns in multidimensional space that
11 cannot be demonstrated by a simple variable value [24].

12 The division of stages contains the potential characteristics of COVID-19. Here, we present
13 three findings. First, at the time of initial diagnosis, the COVID-19 infected patients' physical
14 conditions are similar, regardless of final survival or death. In Figure 4, the distance between
15 stage 1 for the death class and the survival class is small, and the two even overlap. This
16 indicates that outcome prediction is likely not accurate at the time of infection. Second, the
17 physical condition of patients who eventually die changes less than that of those who eventually
18 survive. We conclude this from CH index values, where the CH value of the survival class is
19 larger than that for the death class. Third, mortality rate varies by stage. For example, if the
20 patient is classified into the death class and is at stage 1, there is still hope of survival, as shown
21 by the green triangle sample in Figure 4. However, if the patient is in stage 3 or 4, he or she is
22 very likely to die. Based on estimating the current stage of a patient, doctors will be given a

1 reference, which can help them assess a patient's current situation. Based on that, doctors can
2 carry out targeted treatment and reasonable resource allocation more easily. Thus, the ultimate
3 goal of this study, helping improve the quality of medical care, can be achieved.

4 Meanwhile, by ranking 40 biomarkers according to the degree of correlation with COVID-
5 19 (Table 7), we have found the biomarkers which are more relevant to COVID-19. The top 10
6 are: Lymph, LDH, hs-CRP, Indirect Bilirubin, Creatinine, INR, Serum Sodium, eGFR, Serum
7 Chlorine and Albumin. For each marker, we gave its reference value in each stage, shown in
8 Table 6. Different markers have unique trends in different stages. Combining the correlation
9 analysis with the reference value analysis, we found that the critical COVID-19 patients are
10 usually accompanied by low values of lymph, eGFR, albumin and Serum Sodium, high values
11 of LDH, hs-CRP, indirect bilirubin, creatinine and INR. For example, in the critical stage 4, the
12 average lymph (%) is just 4 and the average LDH (U/l) is up to 499. Besides, there are three
13 major complications of COVID-19 patients - myocardial injury, liver function injury and renal
14 function injury. We got the conclusions separately through the value of 1) LDH, 2) albumin and
15 indirect bilirubin, 3) serum sodium, serum chlorine and creatinine in different stages.

16 Further, there is room for further improvement. First, because of the data limitations, our
17 method may face risk of bias, because data-driven methods are easily influenced by different
18 source of data. For example, the results may vary when using different datasets [13]. Second,
19 we hope to design a more interpretable model, rather than using a complex DL black-box model.
20 Meanwhile, we hope to enlighten the relevant researchers to further study these 4 stages and
21 present more clinical explanations. In particular, we expect to be able to give specific treatments
22 for different stages. Targeted treatment is significant for both patient rehabilitation and the

1 reasonable allocation of medical resources.

2

3 **Conclusions**

4 The sudden outbreak and epidemic of COVID-19 has led to worldwide suffering and shortages
5 of medical resources. In this paper, we propose T-LSTM to predict patient outcomes with high
6 accuracy - 98%, 95% and 93% at 3, 6, and 9 days, which will enable reasonable allocation of
7 medical resources. T-LSTM can effectively model the irregular sampled time series in blood
8 test samples of COVID-19 patients and predict more accurately than existing baselines.
9 Meanwhile, we identified four COVID-19 stages. We ranked biomarkers according to
10 correlations to the outcomes of patients, gave the reference values of top 10 biomarkers for each
11 stage and have found 3 complications. The top 10 biomarkers are: Lymph, LDH, hs-CRP,
12 Indirect Bilirubin, Creatinine, INR, Serum Sodium, eGFR, Serum Chlorine and Albumin. The
13 complications are myocardial injury, liver function injury and renal function injury. By
14 analyzing patients' life conditions at different stages, doctors can choose specific, targeted
15 treatments. Future work will focus more on the treatments in different stages. Meanwhile, more
16 real clinical data are expected to be available for model validation.

17

18 **Abbreviations**

19 COVID-19 (Corona Virus Disease 2019)

20 WHO (World Health Organization)

21 LDH (Lactic Dehydrogenase)

22 hs-CRP (High-sensitivity C-reactive Protein)

- 1 lymph (Lymphocytes)
- 2 DL (Deep Learning)
- 3 RNN (Recurrent Neural Network)
- 4 LSTM (Long Short-Term Unit)
- 5 T-LSTM (Time-aware Long Short-Term Memory)
- 6 PNN (Probabilistic Neural Network)
- 7 RBFNN (Radial Basis Function Neural Network)
- 8 GRNN (Generalized Regression Neural Network)
- 9 BPNN (Back Propagation Neuron Network)
- 10 DT (Decision Tree)
- 11 RF (Random Forest)
- 12 XGBoost (eXtreme Gradient Boosting)
- 13 SVM (Support Vector Machines)
- 14 Cox (Cox's Proportional Hazards Regression)
- 15 LR (Linear Regression)
- 16 NB (Naive Bayes)
- 17 LDA (Linear Discriminant Analysis)
- 18 t-SNE (t-distributed Stochastic Neighbor Embedding)
- 19 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)
- 20 MR (Mortality Rate)
- 21 TD (average Time Distance)
- 22 AUC-ROC (the Area Under the Curve of the Receiver Operating Characteristic)

1 CH (Calinski-Harabaz Index)

2 KL divergence (Kullback-Leibler Divergence)

3

4 **Declarations**

5

6 **Ethics approval and consent to participate**

7 The data is under an MIT license (<https://doi.org/10.5281/zenodo.3758806>).

8

9 **Consent for publication**

10 Not applicable.

11

12 **Availability of data and material**

13 The code implementation is publicly available at <https://github.com/scxhhh/COVID-19>. The
14 data is at https://github.com/HAIRLAB/Pre_Surv_COVID_19.

15

16 **Competing interests**

17 No financial competing interests

18

19 **Funding**

20 This work was supported by the Scientific Research Foundation for the Returned Overseas

21 Chinese Scholars, State Education Ministry and UKRI's Global Challenge Research Fund

22 (ES/P011055/1). The funders had no role in study design, data collection, analysis, the writing

1 of the manuscript, or the decision to submit this article for publication.

2

3 **Authors' contributions**

4 C.S. and S.H. conceptualized the idea. S.H., Z.W. and H.L. initialized, conceived and
5 supervised the project. C.S, S.H. and M.S. collected data and implemented the experiments.

6 C.S, S.H, M.S. and Z.W. drafted the manuscript. All authors provided a critical review of the
7 manuscript and approved the final draft for publication.

8

9 **Acknowledgments**

10 This paper is dedicated to those who want to fight COVID-19.

11

12 **References**

13 1. World Health Organization. Coronavirus Disease 2019 (COVID-19) Situation Report 68,
14 28 March 2020 (2020); [https://www.who.int/docs/default-source/coronaviruse/situation-](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200328-sitrep-68-covid-19.pdf)
15 [reports/20200328-sitrep-68-covid-19.pdf](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200328-sitrep-68-covid-19.pdf)

16 2. World Health Organization. Coronavirus Disease 2019 (COVID-19) Situation Report 147,
17 15 June 2020 (2020); [https://www.who.int/docs/default-source/coronaviruse/situation-](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200615-covid-19-sitrep-147.pdf?sfvrsn=2497a605_2)
18 [reports/20200615-covid-19-sitrep-147.pdf?sfvrsn=2497a605_2](https://www.who.int/docs/default-source/coronaviruse/situation-reports/20200615-covid-19-sitrep-147.pdf?sfvrsn=2497a605_2)

19 3. Huang, C. et al. Clinical features of patients infected with 2019 novel coronavirus in Wuhan,
20 China. *Lancet* 395, 497–506 (2020).

21 4. Chen, N. et al. Epidemiological and clinical characteristics of 99 cases of 2019 novel
22 coronavirus pneumonia in Wuhan, China: a descriptive study. *Lancet* 395, 507–513 (2020).

- 1 5. Yang, X. et al. Clinical course and outcomes of critically ill patients with SARS CoV-2
2 pneumonia in Wuhan, China: a single-centered, retrospective, observational study. *Lancet*
3 *Resp. Med.* 8, 475–481 (2020).
- 4 6. Yan L, Zhang H T, Goncalves J, et al. An interpretable mortality prediction model for
5 COVID-19 patients[J]. *Nature, Machine intelligence*, VOL 2, 2020.
- 6 7. S Hochreiter and J Schmidhuber. Long short-term memory. *Neural Computation*,
7 9(8):1735– 1780.
- 8 8. Baytas I M, Xiao C, Zhang X, et al. Patient Subtyping via Time-Aware LSTM Networks[C]
9 the 23rd ACM SIGKDD International Conference. ACM, 2017.
- 10 9. Laurens V D M, Hinton G. Visualizing Data using t-SNE[J]. *Journal of Machine Learning*
11 *Research*, 2008, 9(2605):2579-2605.
- 12 10. Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu. A Density-Based Algorithm for
13 Discovering Clusters in Large Spatial Databases with Noise. *KDD 1996*: 226-231.
- 14 11. Zhibo Wang, Zhezhi He, Milan Shah, Teng Zhang, Deliang Fan, Wei Zhang:Network-based
15 multi-task learning models for biomarker selection and cancer outcome prediction.
16 *Bioinform.* 36(6): 1814-1822 (2020).
- 17 12. Luchen Liu, Haoran Li, Zhiting Hu, Haoran Shi, Zichang Wang, Jian Tang, Ming Zhang.
18 Learning Hierarchical Representations of Electronic Health Records for Clinical Outcome
19 Prediction. *AMIA 2019*.

- 1 13. Wynants L, Calster B V, Bonten M M J, et al. Prediction models for diagnosis and prognosis
2 of covid-19 infection: Systematic review and critical appraisal[J]. *BMJ (online)*, 2020,
3 369:m1328.
- 4 14. Zhou F, Yu T, Du R, et al. Clinical course and risk factors for mortality of adult inpatients
5 with COVID-19 in Wuhan, China: a retrospective cohort study[J]. *The Lancet*, 2020,
6 395(10229).
- 7 15. Clinical Characteristics of 138 Hospitalized Patients With 2019 Novel Coronavirus–
8 Infected Pneumonia in Wuhan, China[J]. *Jama the Journal of the American Medical*
9 *Association*, 2020.
- 10 16. Yang, X., Yu, Y., Xu, J., Shu, H., Xia, J., Liu, H., et al. (2020) Clinical Course and Outcomes
11 of Critically Ill Patients with SARS-CoV-2 Pneumonia in Wuhan, China: A Single-Centered,
12 Retrospective, Observational Study. *The Lancet Respiratory Medicine*, 1-7.
- 13 17. X.Chu, I.F.Ilyas, S.Krishnan, and J.Wang, Data cleaning: Overview and emerging
14 challenges, in *Proc. Int. Conf. Manage. Data SIGMOD Conf.*, San Francisco, CA, USA,
15 Jun./Jul. 2016, pp. 2201–2206.
- 16 18. Williams R J, Zipser D. A Learning Algorithm for Continually Running Fully Recurrent
17 Neural Networks[J]. *Neural Computation*, 1998, 1(2)
- 18 19. Adam, G., Rampásek, L., Safikhani, Z. et al. Machine learning approaches to drug response
19 prediction: challenges and recent progress. *npj Precis. Onc.* 4, 19 (2020).
- 20 20. Jalali, A., Lonsdale, H., Do, N. et al. Deep Learning for Improved Risk Prediction in
21 Surgical Outcomes. *Sci Rep* 10, 9289 (2020).

- 1 21. Shaoqiang Wang, Shudong Wang, Song Zhang, Fangfang Fan, Gewen He. Research on
2 Recognition of Medical Image Detection Based on Neural Network. IEEE Access 8: 94947-
3 94955 (2020)
- 4 22. Junyuan Shang, Cao Xiao, Tengfei Ma, Hongyan Li, Jimeng Sun. GAMENet: Graph
5 Augmented MEMory Networks for Recommending Medication Combination. AAAI 2019:
6 1126-1133.
- 7 23. Christoph Molnar. Interpretable Machine Learning: A Guide for Making Black Box Models
8 Explainable.
- 9 24. Gibney H M J. Analysis of meal patterns with the use of supervised data mining
10 techniques—artificial neural networks and decision trees[J]. The American Journal of
11 Clinical Nutrition, 2008, 88(6): p.1632-1642.
- 12 25. Emily Czachor. WHO Director Warns COVID-19 Pandemic is 'Speeding Up,' Here for
13 'Long Haul' [N]. Newsweek, News. 6/29/2020. [https://www.newsweek.com/who-director-](https://www.newsweek.com/who-director-warns-covid-19-pandemic-speeding-here-long-haul-1514169)
14 [warns-covid-19-pandemic-speeding-here-long-haul-1514169.](https://www.newsweek.com/who-director-warns-covid-19-pandemic-speeding-here-long-haul-1514169)
- 15 26. Sébastien Farcis. Coronavirus: worries and worries about bed shortage in New Delhi[N].
16 Liberation, Reportage. 6/15/2020.
17 [https://www.liberation.fr/planete/2020/06/15/coronavirus-a-new-delhi-inquietude-et-](https://www.liberation.fr/planete/2020/06/15/coronavirus-a-new-delhi-inquietude-et-desarroi-face-a-la-penurie-de-lits_1791300)
18 [desarroi-face-a-la-penurie-de-lits_1791300.](https://www.liberation.fr/planete/2020/06/15/coronavirus-a-new-delhi-inquietude-et-desarroi-face-a-la-penurie-de-lits_1791300)
- 19 27. Katherine Fung. Arizona Hits Record-High Hospital Capacity As Coronavirus Cases
20 Climb[N]. Newsweek, News. 6/29/2020. [https://www.newsweek.com/arizona-hits-record-](https://www.newsweek.com/arizona-hits-record-high-hospital-capacity-coronavirus-cases-climb-1511578)
21 [high-hospital-capacity-coronavirus-cases-climb-1511578.](https://www.newsweek.com/arizona-hits-record-high-hospital-capacity-coronavirus-cases-climb-1511578)

- 1 28. Sujath R, Chatterjee J M, Hassanien A E. A machine learning forecasting model for COVID-
2 19 pandemic in India[J]. Stochastic Environmental Research and Risk Assessment, 2020(1).
- 3 29. S Dhamodharavadhani, R Rathipriya and Jyotir Moy Chatterjee. COVID-19 Mortality Rate
4 Prediction for India Using Statistical Neural Network Models. Frontiers in Public Health.
5 2020, 8(441).
- 6 30. Panwar H, Gupta P K , Siddiqui M K , et al. Application of deep learning for fast detection
7 of COVID-19 in X-Rays using nCOVnet[J]. Chaos Solitons & Fractals, 2020, 138:109944.
- 8 31. Harsh Panwar, P.K. Gupta, Mohammad Khubeb Siddiqui, Ruben Morales-Menendez,
9 Prakhar Bhardwaj, Vaishnavi Singh. A deep learning and grad-CAM based color
10 visualization approach for fast detection of COVID-19 cases using chest X-ray and CT-Scan
11 images[J]. Chaos Solitons & Fractals, 2020, 140: 110190.
- 12 32. Liang W , Yao J , Chen A , et al. Early triage of critically ill COVID-19 patients using deep
13 learning[J]. Nature Communications. 2020.
- 14 33. R. G. Babukarthik, V. Ananth Krishna Adiga, G. Sambasivam, D. Chandramohan, J.
15 Amudhavel: Prediction of COVID-19 Using Genetic Deep Learning Convolutional Neural
16 Network (GDCNN). IEEE Access, 8: 177647-177666. 2020
- 17 34. Wang L, Wong A. COVID-Net: A tailored deep convolutional neural network design for
18 detection of COVID-19 cases from chest radiography images. arXiv:2003.09871. 2020.
- 19 35. Li L et al (2020) Artificial intelligence distinguishes covid-19 from community acquired
20 pneumonia on chest ct. Radiology, pp 200905, 2020.

- 1 36. de Moraes Batista AF et al. COVID-19 diagnosis prediction in emergency care patients: a
2 machine learning approach. medRxiv. 2020.
- 3 37. Li K et al. The clinical and chest CT features associated with severe and critical COVID-19
4 pneumonia. Investigative radiology, 2020.
- 5 38. Li C et al. Preliminary study to identify severe from moderate cases of COVID-19 using
6 NLR&RDW-SD combination parameter. medRxiv, 2020.
- 7 39. Tang Z et al Severity assessment of coronavirus disease 2019 (COVID-19) using
8 quantitative features from chest CT images. arXiv:2003.11988.
- 9 40. Farid AA, Selim GI, Khater HAA. A novel approach of CT images feature analysis and
10 prediction to screen for corona virus disease (COVID-19). Int J Sci Eng Res 11(3):1–9,
11 2020.
- 12 41. Youssoufa Mohamadou, Aminou Halidou, Pascaline Tiam Kapen: A review of mathematical
13 modeling, artificial intelligence and datasets used in the study, prediction and management
14 of COVID-19. Appl. Intell. 50(11): 3913-3925 (2020).
- 15 42. Kumar R (2020) Accurate prediction of COVID-19 using chest x- ray images through deep
16 feature learning model with smote and machine learning classifiers. In: medRxiv
- 17 43. Corduneanu, C. Integral Equations and Applications[J]. 1991,
18 10.1017/CBO9780511569395.
- 19 44. Chenxi Sun, Shenda Hong, Moxian Song and Hongyan Li. A Review of Deep Learning
20 Methods for Irregularly Sampled Medical Time Series Data. arXiv:2010.12493. 2020.

- 1 45. Radanliev P, Roure D D, Walton R . Data mining and analysis of scientific research data
2 records on covid 19 mortality, immunity, and vaccine development - In the first wave of the
3 Covid-19 pandemic[J]. Diabetes & Metabolic Syndrome: Clinical Research & Reviews,
4 2020.
- 5 46. Gabriel Spadon, Shenda Hong, Bruno Brandoli, Stan Matwin, Jose F. Rodrigues-Jr, Jimeng
6 Sun. "Pay Attention to Evolution: Time Series Forecasting with Deep Graph-Evolution
7 Learning." arXiv preprint arXiv:2008.12833 (2020).
- 8 47. Siuly, S., Zhang, Y. Medical Big Data: Neurological Diseases Diagnosis Through Medical
9 Data Analysis. Data Sci. Eng. 1, 54–64 (2016).

10

11 **Figure Legends**

12 **Figure 1** Examples and statistics of COVID-19 dataset

13 The first block is a line chart of an example in dataset - a 70-year-old female patient. It
14 draws the time series of LHD, lymph and hs-CRP during hospitalization; The second block is
15 the distributions of age, gender, LHD, lymph and hs-CRP of survival class (0) and death class
16 (1); The third block is the statistics about dataset. It contains the counts of time series length,
17 the statistics of overall missing rate and the statistics of each feature's missing rate under
18 different sampling rate.

19 **Figure 2** Structures of the methods

20 The first block shows the structure of RNNs, including basic RNN, LSTM and our T-LSTM;
21 The second block shows the structure that how to use T-LSTM to complete the outcome
22 prediction task (lower grey area) and disease progressing task (upper grey area).

1 **Figure 3** The results of outcome prediction

2 The first line's charts are the AUC-ROC of mortality prediction results using baselines;
3 The second line's chart is the changes of accuracy and loss during training T-LSTM; The third
4 line's charts are the dimension experiments. They show the accuracy of mortality prediction by
5 using different representation dimensions and the effect of representation dimension reduction;
6 The fourth line's charts are the effect when using DBSCAN.

7 **Figure 4** The result of COVID-19 progression

8 This figure shows the four stages of COVID-19 patients by using T-LSTM. The upper
9 clusters are the original clustering of data. The lower are the patient subtyping by using T-LSTM.
10 We can find there are four clusters with distinct boundaries both in death/critical class (red) and
11 survival/general class (blue).

12 **Figure 5** Changes of features in different stages

13 This figure shows the changes of features (Mortality rate, Lymph, LDH, hs-CRP, Indirect
14 Bilirubin, Creatinine, INR, Serum Sodium, eGFR, Serum Chlorine and Albumin) through 4
15 stages. Under different classes, the trends of features are different.

Figures

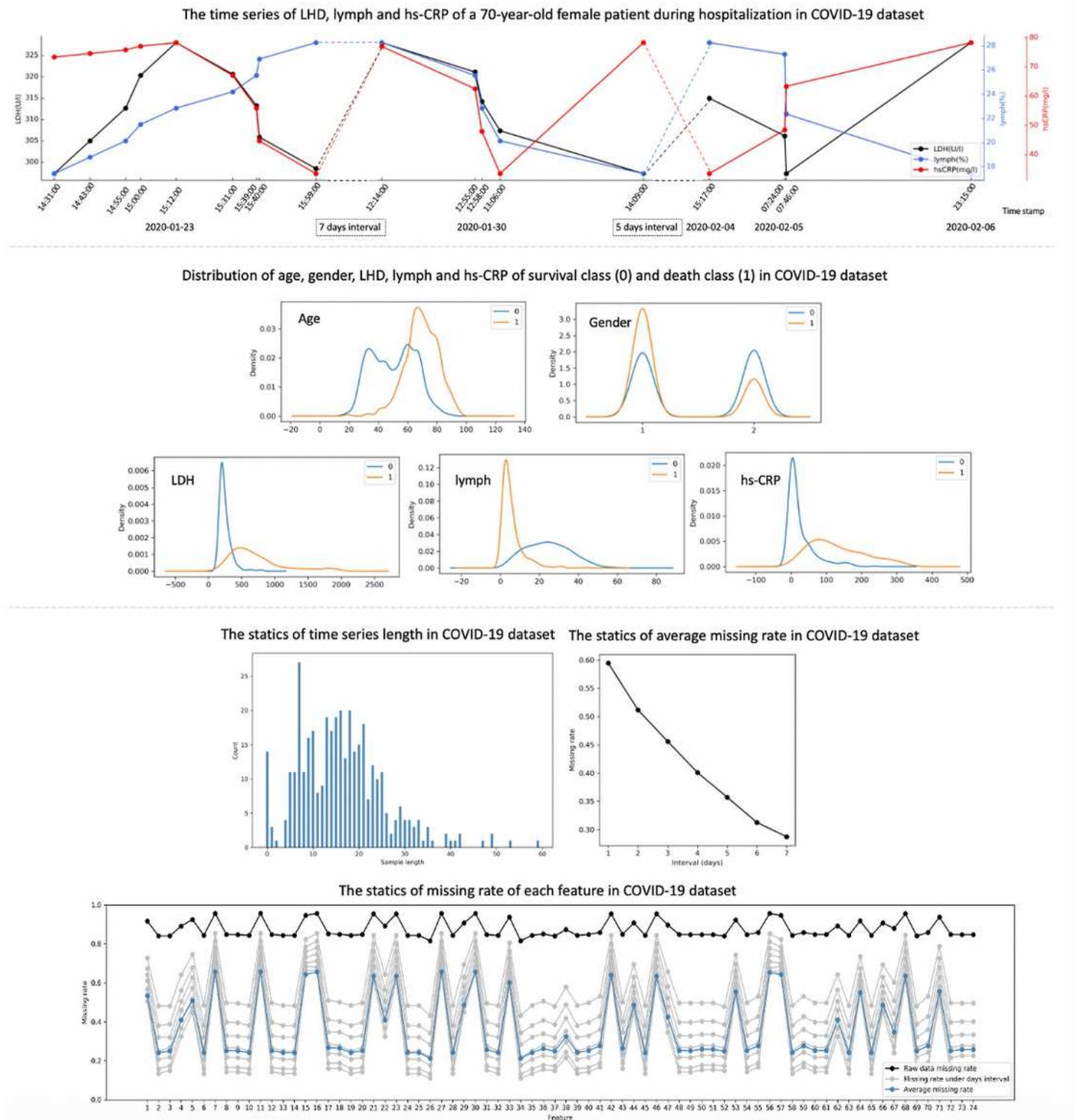
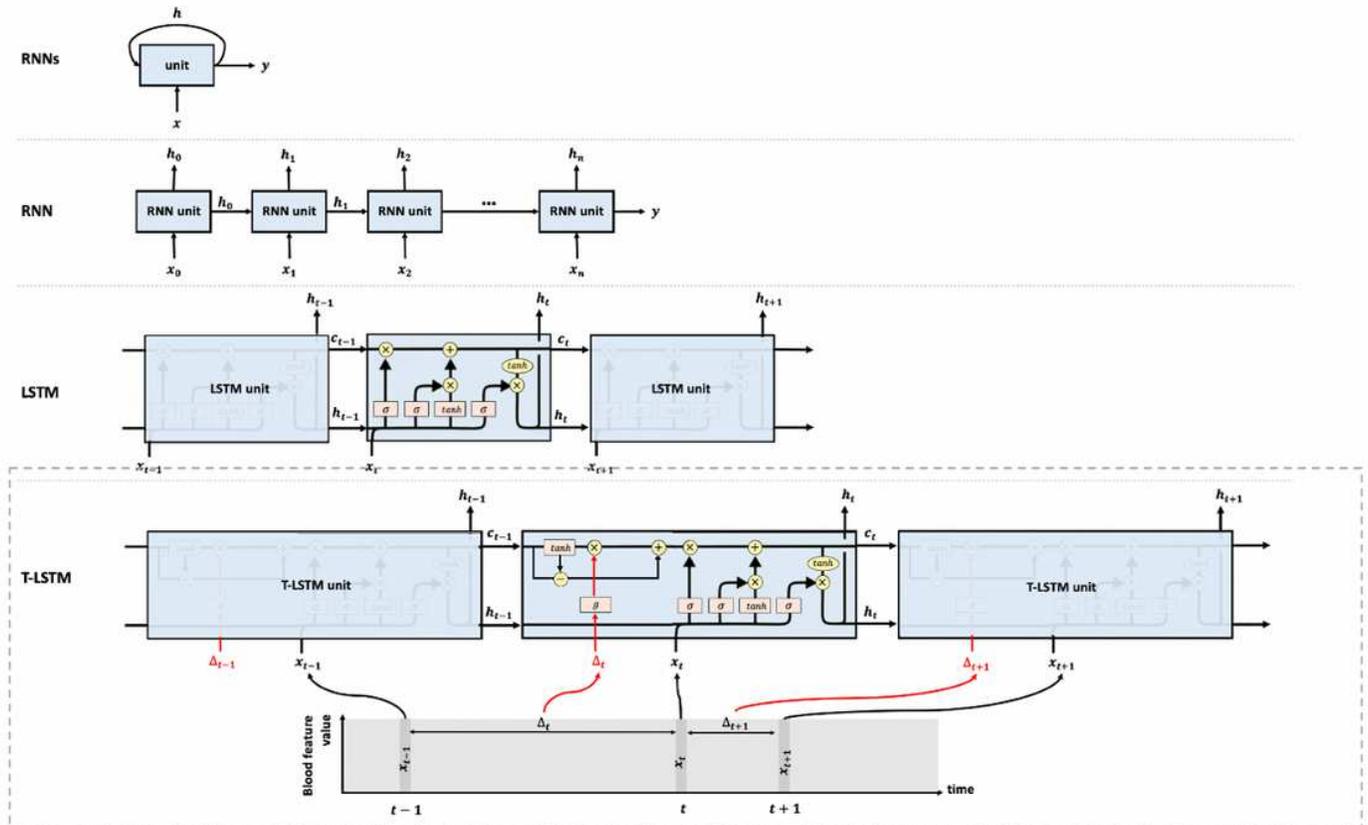


Figure 1

Examples and statistics of COVID-19 dataset The first block is a line chart of an example in dataset - a 70-year-old female patient. It draws the time series of LHD, lymph and hs-CRP during hospitalization; The second block is the distributions of age, gender, LHD, lymph and hs-CRP of survival class (0) and death

class (1); The third block is the statistics about dataset. It contains the counts of time series length, the statistics of overall missing rate and the statistics of each feature's missing rate under different sampling rate.

RNNs structures



Proposed method for predicting COVID-19 progression and patient outcomes

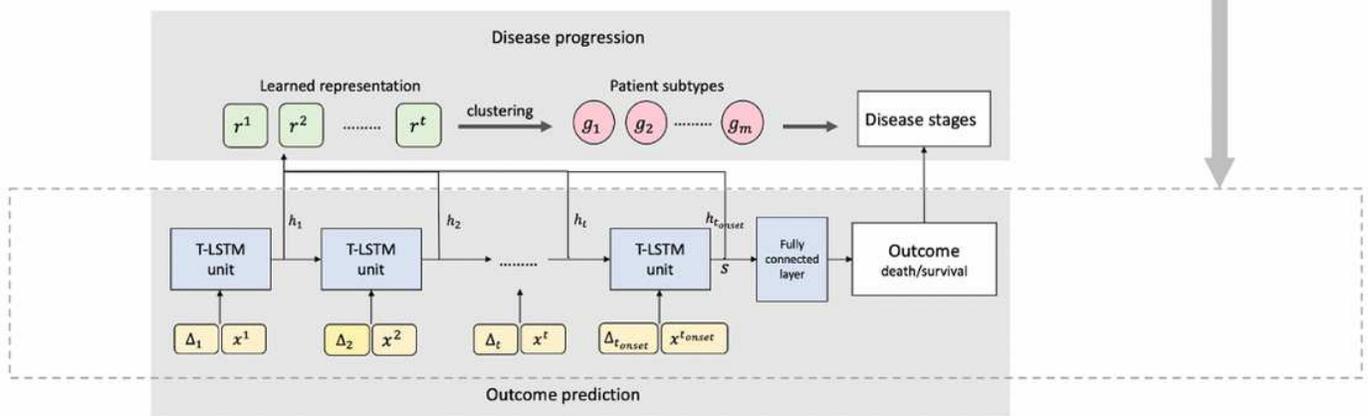
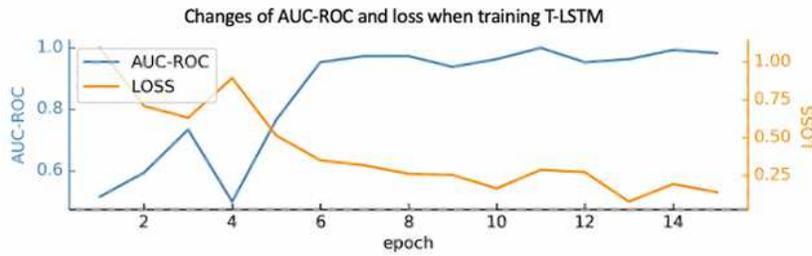
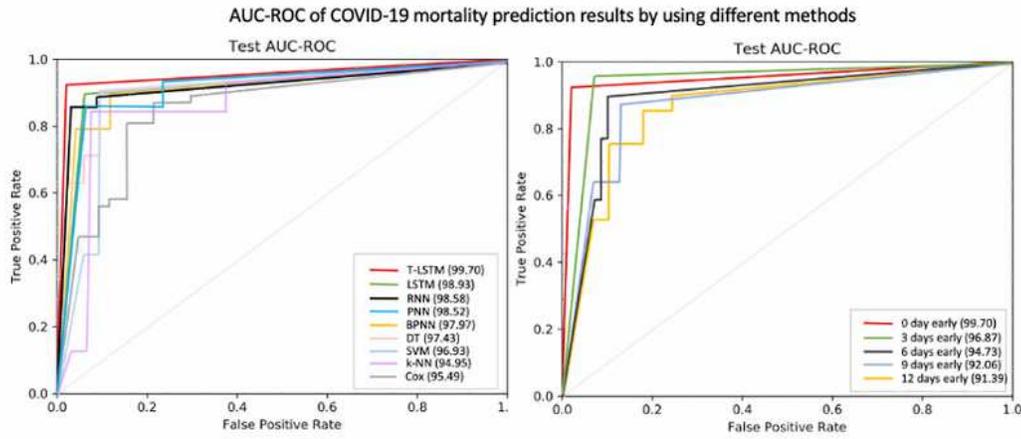
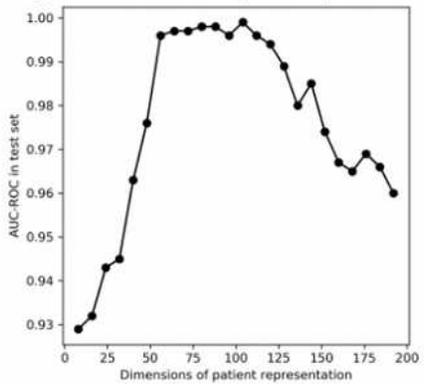


Figure 2

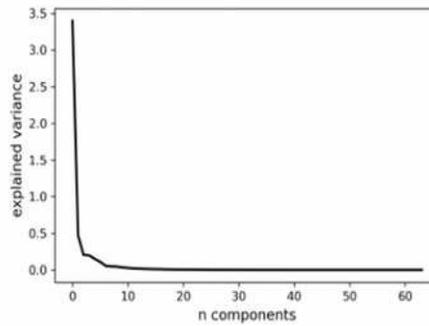
Structures of the methods The first block shows the structure of RNNs, including basic RNN, LSTM and our T-LSTM; The second block shows the structure that how to use T-LSTM to complete the outcome prediction task (lower grey area) and disease progressing task (upper grey area).



AUC-ROC of COVID-19 mortality prediction results by using different dimensions of patient representation



Dimension reduction effect



CH index of DBSCAN clustering effect

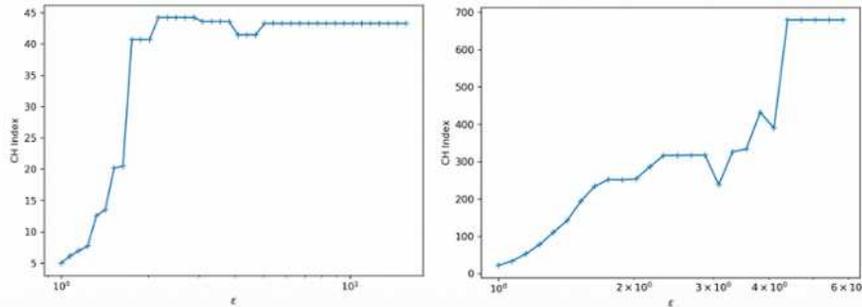
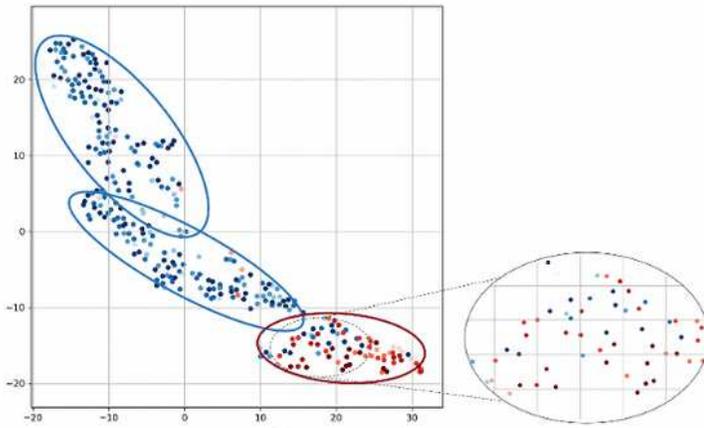


Figure 3

The results of outcome prediction The first line's charts are the AUC-ROC of mortality prediction results using baselines; The second line's chart is the changes of accuracy and loss during training T-LSTM; The third line's charts are the dimension experiments. They show the accuracy of mortality prediction by using different representation dimensions and the effect of representation dimension reduction; The fourth line's charts are the effect when using DBSCAN.

Representation based on raw data



Representation based on method of this paper

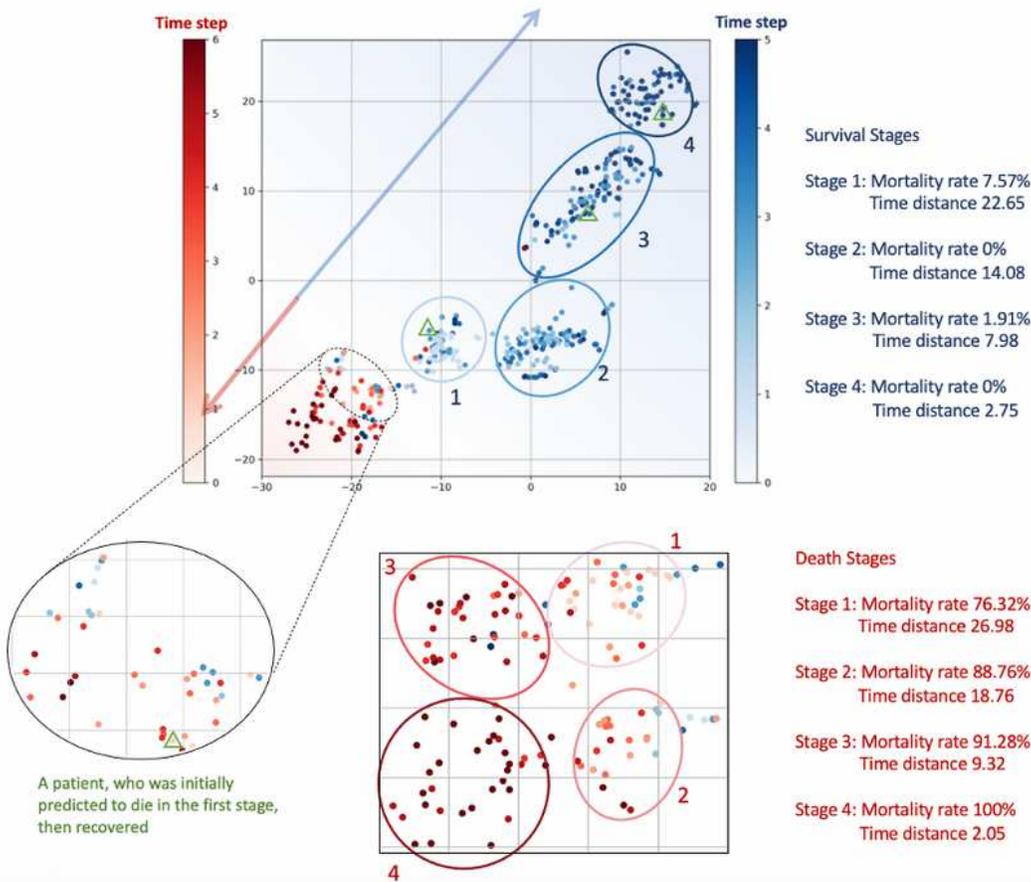
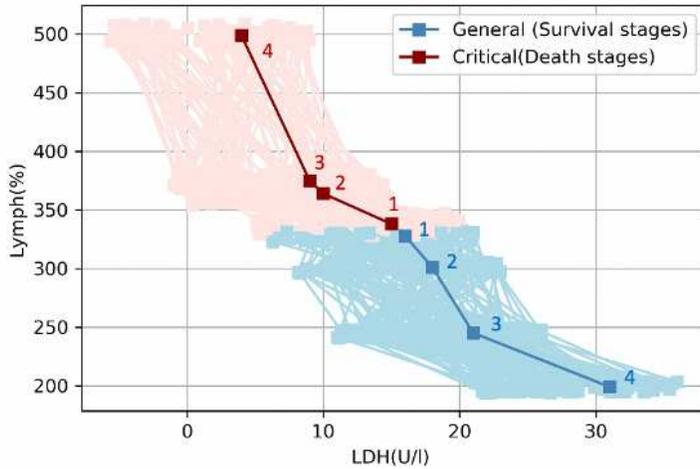


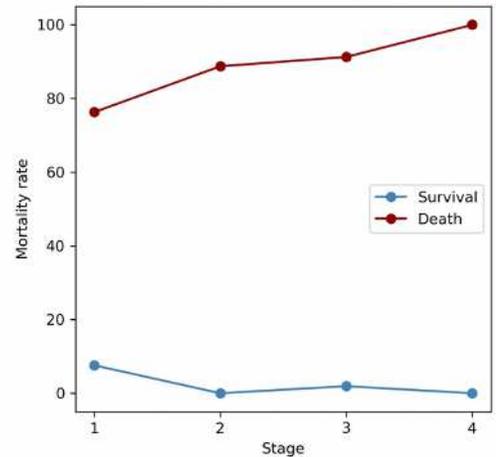
Figure 4

The result of COVID-19 progression This figure shows the four stages of COVID-19 patients by using T-LSTM. The upper clusters are the original clustering of data. The lower are the patient subtyping by using T-LSTM. We can find there are four clusters with distinct boundaries both in death/critical class (red) and survival/general class (blue).

The change of Lymph and LDH of patients in different stages



Mortality rate of patients in different stages



The change of top 10 features of patients in different stages

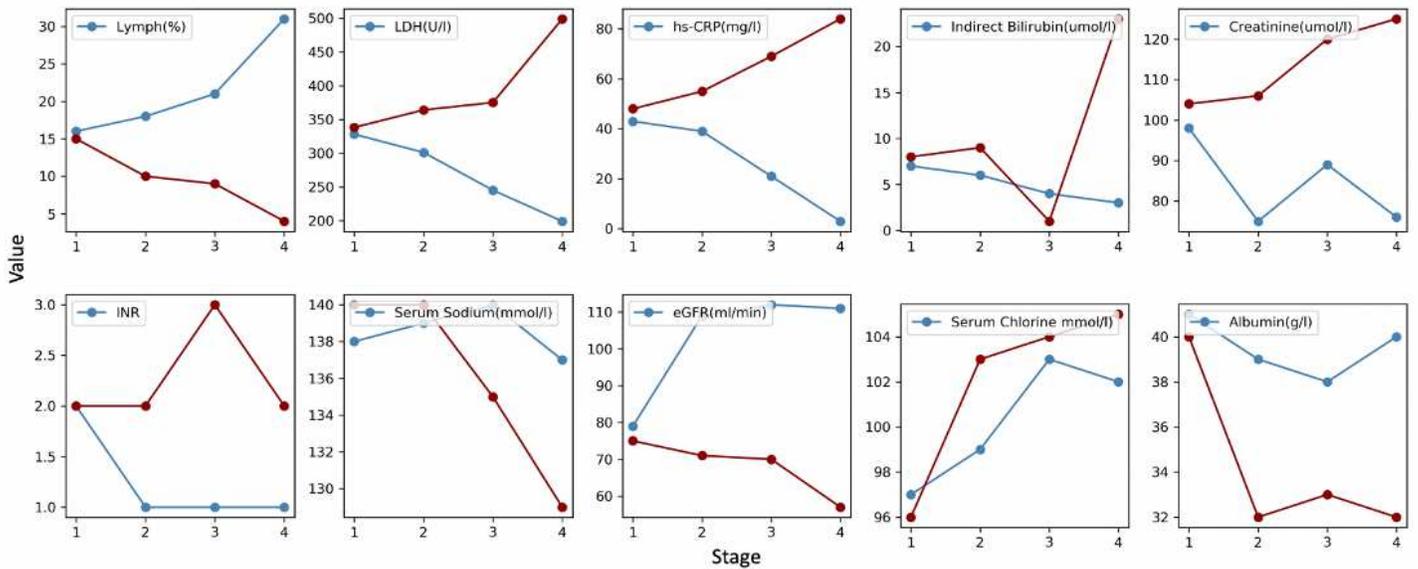


Figure 5

Changes of features in different stages This figure shows the changes of features (Mortality rate, Lymph, LDH, hs-CRP, Indirect Bilirubin, Creatinine, INR, Serum Sodium, eGFR, Serum Chlorine and Albumin) through 4 stages. Under different classes, the trends of features are different.