

iScience, Volume 24

Supplemental information

**A deep learning approach for predicting
severity of COVID-19 patients using
a parsimonious set of laboratory markers**

Vivek Singh, Rishikesan Kamaleswaran, Donald Chalfin, Antonio Buño-Soto, Janika San Roman, Edith Rojas-Kenney, Ross Molinaro, Sabine von Sengbusch, Parsa Hodjat, Dorin Comaniciu, and Ali Kamen

SUPPLEMENTAL MATERIAL

Algorithm 1: Pseudo code for the preprocessing the data.

Input

\mathcal{D}_{demo} patient demographic information
 \mathcal{D}_{labs} clinical laboratory measurements
 $\mathcal{D}_{comorbs}$ co-morbidity information

Output

\mathcal{D}_{all} combined measurements

```
1: procedure PREPROCESS
2:    $mask \leftarrow$  All patients with number of available laboraoty markers  $< 4$ 
3:    $\mathcal{D}_{demo} \leftarrow \mathcal{D}_{demo}[mask]$ 
4:    $\mathcal{D}_{labs} \leftarrow \mathcal{D}_{labs}[mask]$ 
5:    $\mathcal{D}_{comorbs} \leftarrow \mathcal{D}_{comorbs}[mask]$ 
6:    $labs = \{ALT, APTT, AST, Creatinine, CRP, D Dimer, Ferritin, Fibrinogen, Hematocrit,$   
       $INR, LDH, Procalcitonin, Troponin-I, Creatine Kinase, Bilirubin\}$ 
7:   for  $lab \in labs$  do
8:      $\mathcal{D}_{labs}[all, lab] = \log_2(\mathcal{D}_{labs}[all, lab])$ 
9:    $\mathcal{D}_{all} \leftarrow concatenate(\mathcal{D}_{demo}, \mathcal{D}_{labs}, \mathcal{D}_{comorbs})$ 
10:  for  $f \in All\ features\ in\ \mathcal{D}_{all}$  do
11:     $m \leftarrow median(\mathcal{D}_{all}[all, f])$ 
12:     $s \leftarrow interquartile\_range(\mathcal{D}_{all}[all, f])$ 
13:     $\mathcal{D}_{labs}[all, f] = \frac{(\mathcal{D}_{labs}[all, f] - m)}{s}$ 
return  $\mathcal{D}_{all}$ 
```

Algorithm 2: Pseudo code for the training algorithm.

Input

\mathcal{D}_{demo} patient demographic information
 \mathcal{D}_{labs} clinical laboratory measurements
 $\mathcal{D}_{comorbs}$ co-morbidity information
 $\mathcal{D}_{severity}$ severity level

Output

\mathcal{M}_{best} Deep Profiler model

```
1: procedure TRAIN
2:    $epoch \leftarrow 0$ 
3:    $\mathcal{D}_{all}, \mathcal{D}_{mask} \leftarrow \text{PREPROCESS}(\mathcal{D}_{demo}, \mathcal{D}_{labs}, \mathcal{D}_{comorbs})$   $\triangleright$  Pre-process patient information into one container
4:    $train, test \leftarrow \text{SPLIT}(\mathcal{D}_{all}, \mathcal{D}_{severity}, 9)$ 
5:
6:   while  $epoch \leq \text{MaxEpochs}$  do
7:      $\mathcal{O}_{labels}, \mathcal{O}_{recon}, \mathcal{O}_{mean}, \mathcal{O}_{var} \leftarrow \text{DEEPPROFILER}(\mathcal{M}, \mathcal{D}_{all}[train])$ 
8:      $\mathcal{L}_{mse}, \mathcal{L}_{KL} \leftarrow \text{RECONSTRUCTIONLOSS}(\mathcal{O}_{recon}, \mathcal{O}_{mean}, \mathcal{O}_{var}, \mathcal{D}_{all}[train], \mathcal{D}_{mask}[train])$ 
9:      $\mathcal{L} \leftarrow \mathcal{L}_{mae} + \beta \cdot \mathcal{L}_{KL}$ 
10:    if  $epoch > \text{LabelTrainingEpoch}$  then
11:       $\mathcal{L}_{severity} \leftarrow \text{LABELINGLOSS}(\mathcal{O}_{labels}, \mathcal{D}_{severity}[train])$ 
12:       $\mathcal{L} \leftarrow \mathcal{L} + \lambda \cdot \mathcal{L}_{severity}$ 
13:       $\Delta\mathcal{M} \leftarrow \nabla_{\mathcal{M}}\mathcal{L}$   $\triangleright$  Use backpropagation to compute the gradient
14:       $\mathcal{M} \leftarrow \mathcal{M} + \alpha \cdot \Delta\mathcal{M}$   $\triangleright$  Update the model parameters
15:
16:    if  $\text{MODELScore}(\mathcal{M}, \mathcal{D}_{all}[test], \mathcal{D}_{severity}[test]) > \text{MODELScore}(\mathcal{M}_{best}, \mathcal{D}_{all}[test], \mathcal{D}_{severity}[test])$  then
17:       $\mathcal{M}_{best} \leftarrow \mathcal{M}$ 
18:
19:     $epoch \leftarrow epoch + 1$ 
20: return  $\mathcal{M}_{best}$ 
```

Algorithm 3: Pseudo code for computing the training loss function.

Input

\mathcal{O}_{recon} reconstructed data obtained using deep profiler
 \mathcal{O}_{mean} patient latent vectors obtained using deep profiler
 \mathcal{O}_{var} logarithmic of patient latent vector variances obtained using deep profiler
 \mathcal{D}_{all} pre-processed patient data
 \mathcal{D}_{mask} indicator of missing patient measurements

Output

\mathcal{L}_{mae} reconstruction loss
 \mathcal{L}_{KL} regularization loss

```
1: procedure RECONSTRUCTIONLOSS
2:    $N \leftarrow \text{length}(\mathcal{D}_{all})$ 
3:   for  $i \in \mathcal{D}_{mask}$  do ▷ Obtain a indicator array of available measurements
4:      $mask[i, j] \leftarrow 0$  if  $\mathcal{D}_{mask}[patient = i, measurement = j]$  is null else 1
5:      $\mathcal{L}_{mae} \leftarrow \frac{1}{N} \sum |\mathcal{O}_{recon}[mask] - \mathcal{D}_{all}[mask]|$ 
6:      $\mathcal{L}_{KL} \leftarrow -0.5 \cdot \frac{1}{N} \sum (1 + \mathcal{O}_{var} - \|\mathcal{O}_{mean}\|^2 - \exp(\mathcal{O}_{var}))$ 
7:   return  $\mathcal{L}_{mae}, \mathcal{L}_{KL}$ 
```

Input

\mathcal{D}_{demo} patient demographic information
 \mathcal{D}_{labs} clinical laboratory measurements
 $\mathcal{D}_{comorbs}$ co-morbidity information

Output

\mathcal{D}_{all} combined measurements

```
1: procedure LABELINGLOSS
2:    $N \leftarrow \text{length}(\mathcal{D}_{all})$ 
3:    $\mathcal{L} = -\frac{1}{N} \sum (w_i [y_i \cdot \log \sigma(x_i) + (1 - y_i) \cdot \log(1 - \sigma(x_i))])$ 
4:   return  $\mathcal{L}$ 
```

Algorithm 4: Pseudo code for the evaluating the model.

Input

\mathcal{M} Deep Profiler model
 \mathcal{D}_{all} pre-processed patient measurements
 $\mathcal{D}_{severity}$ severity level

Output

\mathcal{S} Deep Profiler model score

```
1: procedure MODELScore
2:    $\mathcal{O}_{labels}, \mathcal{O}_{recon}, \mathcal{O}_{mean}, \mathcal{O}_{var} \leftarrow \text{DEEPPROFILER}(\mathcal{M}, \mathcal{D}_{all})$ 
3:    $\mathcal{S} \leftarrow 0$ 
4:   for  $level \in \{1, 2, 3, 4\}$  do
5:      $\mathcal{S} \leftarrow \mathcal{S} + 0.1 \cdot level \cdot \text{AUROC}(\mathcal{D}_{severity}[level], \mathcal{O}_{labels}[level])$ 
6:   return  $\mathcal{S}$ 
```

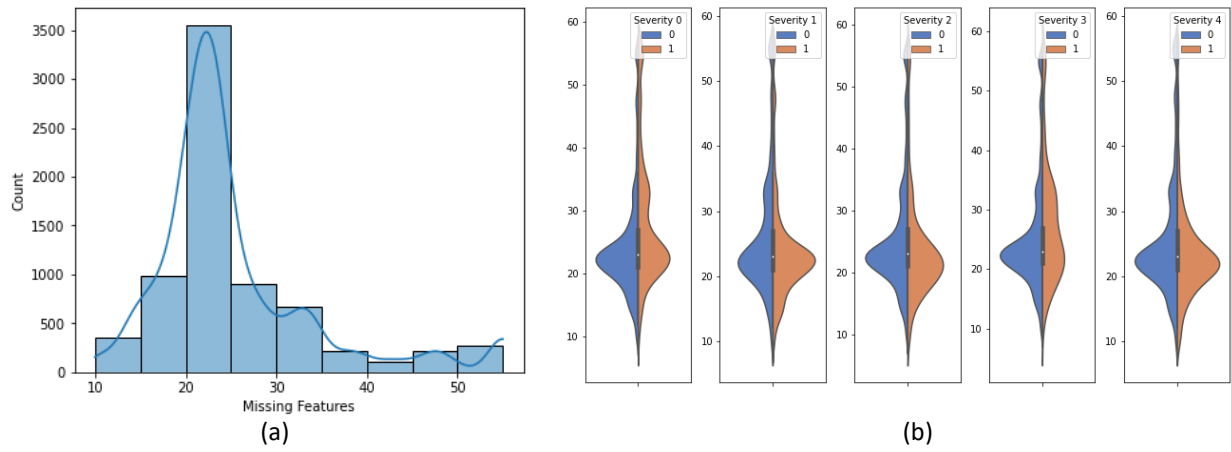


Figure S1: Studying the distribution of missing data over the (a) shows the histogram of the patient count with number of missing input features (NOTE: the patients with less than 4 features are already removed) (b) shows the distribution of the patients with missing input features by different severity levels.

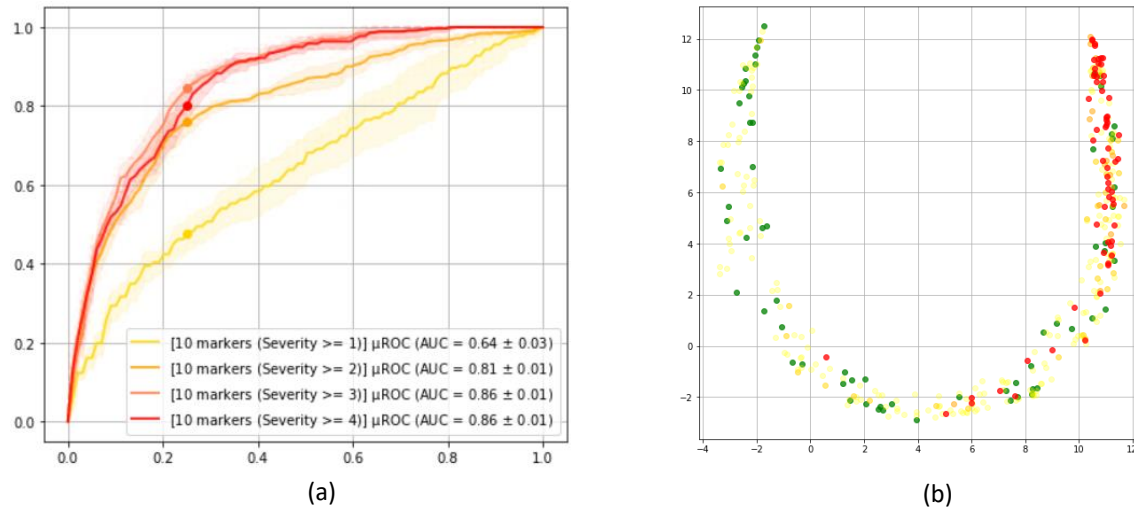
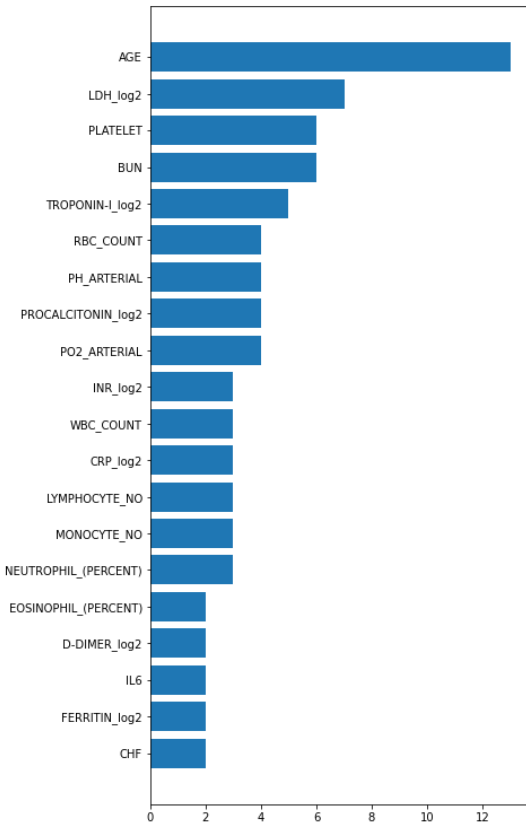
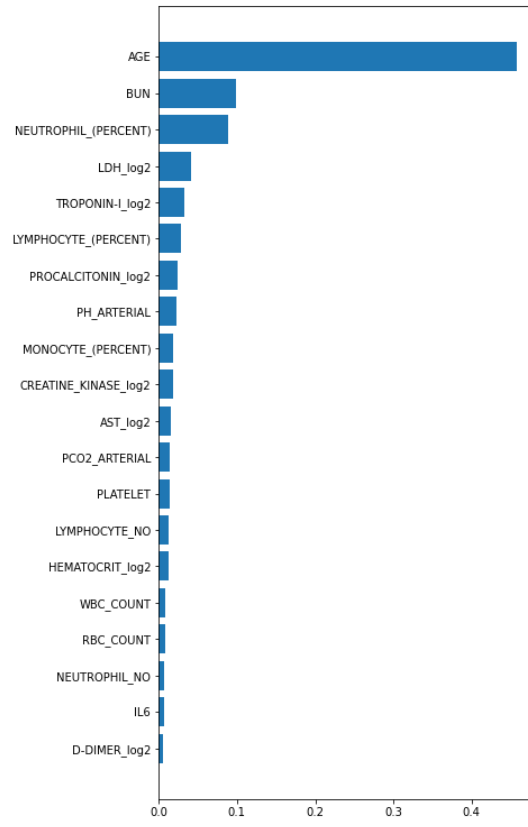


Figure S2: Studying the impact of the missing data on the model (a) shows the performance of the 10-marker deep profiler on a subset of 300 patients in the testing data for which all the 10 measurements were available within first 24 hours of admission. Out of 300 patients, number of patients with severity ≥ 1 , ≥ 2 , ≥ 3 , ≥ 4 is 244, 64, 45 and 36 respectively (b) shows the UMAP visualization of the same 300 patients on the deep profiler latent space for first model in the ensemble.



(a)



(b)

Figure S3: Feature importance of the machine learning models (a) XGBoost (b) Random Forest Regression.