1  **The Selection of Indicators from Initial Blood Routine Test Results to Improve the**

2  **Accuracy of Early Prediction of COVID-19 Severity**

3  **Running Title:** COVID-19 Severity Prediction using MCDM Algorithm

4  Jiaqing Luo[1#], Lingyun Zhou[3#*], Yunyu Feng[4], Bo Li[5], Shujin Guo[2*]

5  1. School of Computer Science and Engineering, University of Electronic Science and

6  Technology of China, Chengdu 611731, China.

7  2. The Geriatric Respiratory department, Sichuan Provincial People's Hospital, University of

8  Electronic Science and Technology of China, Chengdu 611731, China.

9  3. Center of Infectious Diseases, West China Hospital of Sichuan University. Chengdu 610041,

10  China.

11  4. State Key Laboratory of Biotherapy and Cancer Center, West China Hospital, Sichuan

12  University and Collaborative Innovation Center. Chengdu 610041, China

13  5. Department of Otorhinolaryngology, Head & Neck Surgery, West China Hospital, Sichuan

14  University, Chengdu, 610041, China

15  # Jiaqing Luo and Lingyun Zhou contributed equally to this work and are joint first authors.

16  **Address for correspondence:**

17  * Lingyun Zhou, Center of Infectious Diseases, West China Hospital of Sichuan University.

18  Chengdu 610041, China.

19  E-mail: 4423925@qq.com

20  * Shujin Guo, The Geriatric Respiratory department, Sichuan Provincial People's Hospital,

21  University of Electronic Science and Technology of China, Chengdu 611731, China.

22  E-mail: shujinguo@126.com

23

24 **Abstract**

25  Early prediction of disease severity is important for effective treatment of COVID-19. We

26 determined that age is a key indicator for severity predicting of COVID-19, with an accuracy of

27 0.77 and an AUC of 0.92. In order to improve the accuracy of prediction, we proposed a Multi

28 Criteria Decision Making (MCDM) algorithm, which combines the Technique for Order of

29 Preference by Similarity to Ideal Solution (TOPSIS) and Naïve Bayes (NB) classifier, to further

30 select effective indicators from patients' initial blood test results. The MCDM algorithm selected

31 3 dominant feature subsets {Age, WBC, LYMC, NEUT}, {Age, WBC, LMYC} and {Age,

32 NEUT, LYMC}. Using these feature subsets, the optimized prediction model could achieve an

33 accuracy of 0.82 and an AUC of 0.93. This result indicated that using age and the indicators

34 selected by the MCDM algorithm from blood routine test results can effectively predict the

35 severity of COVID-19 at an early stage.

36 **Keywords:** Coronavirus disease 2019 (COVID-19), Severity, Blood routine test, Multiple

37 Criteria Decision Making (MCDM).

38

39 **Introduction**

40 Currently, more than 40 million people worldwide are infected with the SARS-Cov-2 virus, and

41 more than 10 million people are suffering from Coronavirus disease 2019 (COVID-19) and are

42 receiving treatment. This poses a huge threat to the health and lives of people all over the world,

43 and brings unprecedented pressure to the medical system. Many infected patients cannot receive

44 timely and effective treatment, and it will also reduce the treatment efficiency of other

45 emergency patients.

46 Patients with suspicious symptoms and epidemiological history first visit the fever clinic of the

47 community hospital. They usually undergo 3 initial tests: SARS-Cov-2 RNA confirms SARS-

48 Cov-2 infection (*1*), blood routine test and chest CT scan to initially assess the severity of

49 COVID-19 (*2-4*). The timely and effective triage of COVID-19 patients based on the results of

50 the 3 initial tests is of great significance for maintaining emergency capacity and optimizing

51 treatment plans.

52 Although most COVID-19 patients are Mild-Moderate cases and can recover on their own, about

53 14% of patients are Severe cases, and 5% of patients are Critically Severe cases (*5*). Severe-

54 Critically Severe cases usually develop Acute Respiratory Distress Syndrome (ARDS) or

55 Multiple Organ Dysfunction Syndrome (MODS) within 2 weeks of infection (*6*), which

56 consumes most of medical resources and leads to a high case fatality rate (up to 49%) (*5*). Early

57 prediction of the severity of COVID-19 can not only help quickly triage patients (i.e., quarantine,

58 hospital admission or ICU assignment, etc.), but also optimize the use of medical resources and

59 timely medical intervention. Previous studies have used multiple indicators to predict the

60 severity of COVID-19 (i.e., older age, pulmonary micro-thrombosis, increased inflammatory

61 factors (C-reactive protein (CRP), IL-6), hyper-lactic acidemia, D-dimer progressive heightened,

62    decreased lymphocyte count (especially CD8+ T cell count) and short-term progression of lung

63    lesions, etc.). However, the collection of these indicators requires multiple tests and takes a lot of

64    time.

65    Of all the initial tests, blood routine test is the worldwide common test, and the results are

66    usually available within 2 hours. In this paper, we tried to select features from blood test results

67    to predict the severity of COVID-19 quickly and accurately. Specifically, we first defined feature

68    selection as a Multiple Criteria Decision Making (MCDM) problem that considers the correlation

69    between input features, and the correlation between input and output features, and then combined

70    the Technique for Order of Preference by Similarity to Ideal Solution (TOPSIS) and Naïve Bayes

71    (NB) classifier to achieve the highest prediction accuracy with the fewest features.

72    Our early prediction of the severity of COVID-19 based on the clinic characteristics of patients

73    can improve the efficiency and accuracy of emergency triage of patients, thereby effectively

74    supplementing and improving the overall management of COVID-19.

75

76 **Methods**

77 *Patient enrollment and study design*

78 This retrospective study was approved by the ethics committee of Sichuan Provincial People's

79 Hospital. We collected 196 COVID-19 patients diagnosed according to WHO guidance (*7*) in

80 Wuhan Red Cross Hospital from February 1, 2020 to March 15, 2020. Written or oral informed

81 consent was obtained from patients.

82 *Definitions*

83 COVID-19 was confirmed by detecting SARS-CoV-2 RNA test. According to the 5th edition of

84 the China Guidelines for the Diagnosis and Treatment Plan of COVID-19 Infection by the

85 National Health Commission (Trial Version 5) (*8*), the cases were classified into Mild-Moderate

86 and Severe- Critically Severe.

87 *Data collection*

88 The following information was extracted from each patient: Gender, Age and patients' initial

89 blood routine test results including White Blood Cell Count (WBC), Lymphocyte Count

90 (LYMC), Lymphocyte Ratio (LYMPH), Neutrophil Count (NEUT), Neutrophil Ratio (NEU) and

91 Neutrophil to Lymphocyte Ratio (NLR). The dataset contained 8 input features {Gender, Age,

92 WBC, LYMC, LYMPH, NEUT, NEU, NLR}, and 1 output feature (Severity).

93 *Statistical Analysis*

94 Quantitative variables were expressed as the mean ± standard deviation or the median with

95 interquartile ranges, while categorical variables were expressed as absolute and relative

96 frequencies. The t test or Wilcoxon-test was performed to calculate differences between

97 quantitative data; and $\chi 2$ test was performed to calculate differences between qualitative data.

98 According to the data characteristics, the correlation between clinic characteristics and COVID-

99   19 severity was calculated according to Kendall correlation coefficient (Gender-severity) or

100  Spearman correlation coefficient. Logistic regression analysis was performed for independent

101  variables with collinearity. Wald test was used to determine the joint significance of variables.

102  The standard deviation was used to measure dispersion degrees. Statistical procedures were

103  performed with R statistical software. P values of ≤0.05 were considered significant.

104  ***The MCDM algorithm design and implementation***

105  The proposed algorithm is basically designed for predicting COVID-19 severity, either Mild-

106  Moderate or Severe-Critically Severe case. It leads to reducing computation time, improving

107  prediction performance, and a better understanding of the data in machine learning.

108  It consists of 4 major stages: preprocessing, feature ranking, feature selection and performance

109  evaluation. Preprocessing is the process to refine the collected raw data to de-noise it. Feature

110  ranking is the process of ordering the features by the value of some scoring function, which

111  usually measures feature-relevance. Feature selection aims to choosing a small subset of the

112  relevant features from the original features by removing irrelevant, redundant, or noisy features

113  (*9*). Performance evaluation is to measure the performance of the binary classification by

114  statistical measures, i.e., Accuracy (ACC), True Positive Rate (TPR), False Positive Rate (FPR)

115  and F1 score.

116  ● Preprocessing -We use stratified random sampling to divide the dataset into 2 subsets:

117  training set (80%) and test set (20%). In these 4 stages, we only used the test set for

118  performance evaluation. Suppose there are m input features and n output features. Let

119  $X=\{x|1\le x\le m\}$ be the input feature set and $Y=\{y|m+1\le y\le m+n\}$ be the output feature set.

120  Elements x and y are indexes of features. The feature set is $F=X\cup Y=\{i|1\le i\le m+n\}$. We

121  calculated and visualized a $(m+n)\times(m+n)$ correlation matrix R and a $(m+n)\times(m+n)$ p-

122     value matrix P to show the correlations between all different feature pairs. To simplify the

123     analysis, we then preprocess R in 2 steps. STEP1: We ignored the sign of R[i,j]. Let R[i,j] = |

124     R[i,j]| so that the range of R[i,j] changes from [-1,1] to [0,1], where $i,j \in F$. STEP2: We

125     filtered R through P. For $x \in X$ and $y \in Y$, if P[x,y] = P[y,x] > 0.05, R[x,y] and R[y,x] are

126     not significant. We set R[x,y] = R[y,x] = 0 and R[x,i] = R[i,x] = 1 for $i \in X$.

127 ●  Feature Ranking-We defined a labeled feature set L and initialized with $L = \varnothing$. We iterated

128     the procedure of ranking input features $x \in X$ and moved the first in each ranking from X to

129     L. The ranking criteria includes 2 evaluations: EVAL1: The correlation between input

130     feature $x \in X$ and output feature $y \in Y$, R[x,y] or R[y,x]. EVAL2: The correlation between

131     input feature $x \in X$ and labeled feature $v \in L$, R[x,v] or R[v,x]. This explicitly evaluates

132     multiple conflicting criteria in decision making. We proposed an algorithm to solve this

133     Multiple Criteria Decision Making (MCDM) problem by using the Technique for Order of

134     Preference by Similarity to Ideal Solution (TOPSIS) (*10*), which is a compensatory

135     aggregation method. The algorithm, called MCDM, creates an evaluation matrix E

136     consisting of p criteria and q alternatives, to rank input features. According to Pareto's

137     principle (*11*), the algorithm divide x into the following 2 types:

138     TYPE1: If $|X| > min \{m-1, \lceil 0.8 \times m \rceil\}$, x to be labeled are core features (the top 20%),

139     which should have the lowest R[v,x] from EVAL2, and the highest R[y,x] from EVAL1.

140     The algorithm sorts the elements of sets $L \cup Y$ and X in ascending order to get sequences

141     $\left(r_i\right)_{i=1}^{|L|+n}$ and $\left(c_j\right)_{j=1}^{¿X \vee ¿¿}$, respectively. Let $p = ¿L \vee +n$ and $q = ¿X \vee ¿$, the algorithm extracts a

142     $p \times q$ submatrix E from R such that $E[i,j] = R[r_i, c_j]$. The worst condition of $E[i,j]$ is $w_i = ¿$,

143     and the best condition of $E[i,j]$ is $b_i = ¿$.

144    TYPE2: If $|X| \leq min\{\vee m-1, \lceil 0.8 \times m \rceil\}$, x to be labeled are auxiliary features (the rest

145    80%), which only need to have the lowest R[v,x] from EVAL2. The algorithm sorts the

146    elements of sets $L$ and X in ascending order to get sequences $(r_i)_{i=1}^{|L|}$ and $(c_j)_{j=1}^{¿X\vee¿¿}$,

147    respectively. Let $p=¿L\vee¿$ and $q=¿X\vee¿$, E is a $p \times q$ matrix with $E[i,j]=R[r_i,c_j]$.

148    The algorithm calculates the L2-distance between the target alternative j and the worst

149    condition:

150
$$d_{wj}=\sqrt{\sum_{i=1}^{p}(E[i,j]-w_i)^2} \quad \text{Eq.1}$$

151    It then calculates the distance between j's condition and the best condition:

152
$$d_{bj}=\sqrt{\sum_{i=1}^{p}(E[i,j]-b_i)^2} \quad \text{Eq.2}$$

153    After that, it calculates the similarity to the worst condition:

154
$$s_j=\frac{d_{wj}}{d_{wj}+d_{bj}}, 0 \leq s_j \leq 1 \quad \text{Eq.3}$$

155    $s_j=1$ if and only if alternative j has the best condition, and $s_j=0$ if and only if alternative j

156    has the worst condition. Let $j^¿=arg\underset{j}{max}\{s_j\}$, then $X=X¿c_{j^¿}\}$ and $L=L\cup\setminus\{c_{j^¿}\}$.

157        The pseudocode of the MCDM algorithm is as follows:

| Algorithm MCDM is |
| --- |
| Input: correlation matrix R, number of input features m, number of input features n, input feature set X, output feature set Y |
| Output: labeled feature set L |
| initialize $L=\varnothing$ |
| while $X \neq \varnothing$ do |

```
if |X| > min {m−1, ⌈0.8 × m⌉}

    (r_i)_{i=1}^{|L|+n} ← sort L∪Y and X in ascending order

else

    (r_i)_{i=1}^{|L|} ← sort L in ascending order

    (c_j)_{j=1}^{|X∨i|} ← sort X in ascending order

    extract E from R such that E[i, j] ← R[r_i, c_j]

    for j = 1 to q do  // q is the number of columns of E

        d_{wj} ← Eq.2

        d_{bj} ← Eq.3

        s_j ← Eq.4

        j* ← arg max {s_j}
               j

        X ← X ↓ c_{j*}}

        L ← L∪\{c_{j*}}

    print L

return L
```

165      analyze the performance of feature selection from Accuracy (ACC), True Positive Rate

166      (TPR), False Positive Rate (FPR) and F1 score.

167    ***Evaluation of the predictive value of selected features***

168    According to stratified random sampling, we divided the data set into 2 subsets: 80% of the

169    "training set" and 20% of the "testing set". We used Receiver Operating Characteristic (ROC)

170    curve analysis to calculate the Area Under the Curve (AUC) and use "ROC" package in R to

171    evaluate the prediction accuracy of our model.

172

173    **Results**

174    *Baseline characteristics*

175    We analyzed the data of 196 COVID-19 patients, of which 67 and 129 were male and female

176    patients. After clearing the data set, there is no abnormal data (S-Figure 1). Table 1 lists the

177    detailed baseline characteristics. The mean age of patients was 57.74±15.87 years old. The

178    COVID-19 patients' initial blood routine test results showed that the WBC was 6.75±3.49◊109/

179    L; LYMC was 1.12±0.58◊109/L; LYMPH was 19.91±11.52%; NEUT was 5.13±3.46◊109/L;

180    NEU was 71.34±15.24%; the NLR was 7.45±13.08.

181    *Difference in Age and initial blood test results between Mild-Moderate and Severe-Critically*

182    *Severe groups*

183    According to the 5th edition of the China Guidelines for the Diagnosis and Treatment Plan of

184    COVID-19 Infection by the National Health Commission, we divided patients into 2 groups: 67

185    cases in the Mild-Moderate group, and 129 cases in the Severe-Critically Severe group (Table 1).

186    Comparing Mild-Moderate and Severe-Critically Sever groups, the basal features showed no

187    differences in Gender (p=0.26) (Figure1A). The Severe-Critically Severe group was significantly

188    older than the Mild-Moderate group (p <0.001) (Figure 1B). The initial blood routine test seems

189    to be important for predicting the severity of COVID-19: The Severe-Critically Severe group had

190    a higher WBC level (p=0.02) (Figure1C). The Severe-Critically Severe group had extremely low

191    LYMC (p < 0.001) and LYMPH (p < 0.001) (Figure1D, E). In contrast, NEUT (p < 0.001) and

192    NEU (p < 0.001) in the Severe-Critically Severe group were extremely high (Figure1F, G). As a

193    result, the Severe-Critically Severe group had a higher NLR (p < 0.001) (Figure1H).

194    *Predictive value of age and initial blood test results for COVID-19 severity*

195   By calculating the correlation between clinic characteristics and severity of COVID-19, we

196   found that Age ($r=0.73$, $p=0.01$), WBC ($r=0.24$, $p < 0.01$), NEUT ($r=0.34$, $p < 0.01$), NLR

197   ($r=0.31$, $p < 0.01$) were significantly positively correlated with the severity of COVID-19, while

198   LYMC ($r=-0.55$, $p=0.01$) was significantly negatively correlated with the severity of COVID-19

199   (Figure 2A, B). These results indicated that Age and initial blood routine test results-WBC,

200   LYMC, NEUT, NLR, might be important for predicting the severity of COVID-19.

201   Wald test showed that only Age was the key indicator in predicting the severity of COVID-19

202   (Table2). Using stratified random sampling, we generated the Receiver Operating Characteristic

203   (ROC) curve to evaluate the predictive values: 80% for the "training set" and 20% for the

204   "testing set". Using {Age} for prediction, we can obtain an accuracy of 0.77, and an Area Under

205   the Curve (AUC) of 0.92 (Figure2C). Through dispersion analysis, we found that WBC, LYMC

206   and LYMPH may be able to optimize prediction performance (Table3, Table4). The ROC curve

207   showed that {Age, WBC, LYMC} had an accuracy of 0.82 and an AUC of 0.93 (Figure2D).

208   ***Details of the MCDM algorithm to predict the severity of COVID-19***

209   The MCDM algorithm and Logistic regression analysis have obtained consistent results: Age

210   was a key indicator in predicting the severity of COVID-19. In addition, the MCDM algorithm

211   verified that the {Age, WBC, LYMC} subset is one of the index sets with the highest prediction

212   accuracy.

213   Preprocessing (Figure3A) - In the COVID-19 data set, $m=8$ and $n=1$. The $9 \times 9$ correlation

214   matrix R, The $9 \times 9$ p-value matrix P and the range of R[i,j] for $i, j \in F$ becomes [0,1]. Since

215   P[1,9]=P[9,1]=0.1442>0.05 , R[1,9] and R[9,1] are not significant, R[1,9]=R[9,1]=0,

216   R[1,1:8]=ones(1,8) and R[1:8,1]=ones(8,1).

217     Feature Ranking (Figure3B) - When $|X|=8>min\{8-1,\lceil 0.8\times 8\rceil\}=7$, $L\cup Y=\varnothing\cup\{9\}=\{9\}$ and

218     $X=\{1,\ldots,8\}$. Then, we have, $(r_i)_{i=1}^1=(9)$ and $(c_j)_{j=1}^8=(1,\ldots,8)$. Since $p=|L|+n=1$ and

219     $q=|X|=8$, E is a $1\times 8$ submatrix of R. When $|X|=5<7$, $L=\{2,3,4\}$ and $X=\{1,5,6,7,8\}$. Then,

220     we have $(r_i)_{i=1}^3=(2,3,4)$ and $(c_j)_{j=1}^5=(1,5,6,7,8)$. Since $p=|L|=3$ and $q=|X|=5$, E is a $3\times 5$

221     submatrix of R. When $|X|=8>7$ , $w_i=1$ and $b_j=0$. By Eq. 1 and Eq.2, we calculated

222     $d_{w2}=0.5913$ and $d_{b2}=0.4087$. By Eq. 3, we have $s_2=0.5913$. When $|X|=5<7$, $w_i=1$ and $b_j=0$.

223     By Eq.1 and Eq.2, we calculated $d_{w6}=1.1871$ and $d_{b6}=0.9912$. By Eq. 3, we got $s_6=0.5450$.

224     Feature Selection (Figure3C) - When 4 features {2,5,8,4} are selected, the accuracy of EVAL1

225     reached a peak of 0.803. Interestingly, with less features {2,3,4}, the accuracy of

226     EVAL1+EVAL2 can reach a higher 0.815.

227     Performance Evaluation (Figure3D) - {2,3,4} has the lowest number of features, but the highest

228     score among multiple performance metrics. We can see that the accuracy of {2,5,8,4,7,6,3},

229     {2,5,8,4} and {2,3,4} are 0.74, 0.82 and 0.87, respectively. We can also see that the F1 score of

230     {2,5,8,4,7,6,3}, {2,5,8,4} and {2,3,4} are 0.67, 0.72 and 0.78, respectively.

231     ***Influence of dataset uncertainty on the feature selection of the MCDM algorithm***

232     To test the stability of the algorithm and observe the influence of the dataset uncertainty on

233     feature selection, we divided the data set 100 times (80% training set and 20% test set) and

234     repeatedly run the algorithm. The average number of features selected by 3 different criteria,

235     EVAL1, EVAL1 (subset) and EVAL1+EVAL2 (subset) are 6.58 (95% CI: 6.48 - 6.68), 3.26

236     (95% CI: 3.01 - 3.51) and 3.52 (95% CI: 3.40 - 3.64), respectively (Figure4A). The criteria,

237     EVAL1+EVAL2 (subset), adopted by the MCDM algorithm improved most performance

238     metrics. The metrics (ACC, TPR, FPR and F1 score) of EVAL1+EVAL2 (subset) are 0.81 (95%

239    CI: 0.80 - 0.82), 0.69 (95% CI: 0.67 - 0.71), 0.09 (95% CI: 0.08 - 0.11) and 0.75 (95% CI: 0.73 -

240    0.77) respectively, while those of EVAL1 are 0.75 (95% CI: 0.74 - 0.77), 0.60 (95% CI: 0.58 -

241    0.62), 0.07 (95% CI: 0.06 - 0.09) and 0.71(95% CI: 0.70 - 0.73) respectively (Figure4B).

242    Although dataset uncertainties have an influence on feature selection, there were still 3 subsets:

243    {Age, WBC, LYMC, NEUT} with a selection rate of 44%, {Age, NEUT, LYMC} with a

244    selection rate of 38%, and {Age, WBC, LYMC} with a selection rate of 9%, which dominated

245    EVAL1+EVAL2 (subset) feature selection. These 3 subsets can achieve high accuracy with a

246    small number of features (Figure4C).

247    ***Predictive value of the features selected by the MCDM algorithm***

248    Using stratified random sampling, we generated ROC curves to evaluate the predictive values of

249    the subsets selected by the MCDM algorithm: 80% for the "training set" and 20% for the "testing

250    set". Our analysis results showed that {Age, WBC, LYMC, NEUT} (Figure5A), {Age, NEUT,

251    LYMC} (Figure5B) and {Age, WBC, LYMC} (Figure5C) all achieved 0.82 accuracy and 0.93

252    AUC. The MCDM algorithm can steadily and accurately select Age and other features from

253    initial blood routine test results to predict the severity of COVID-19.

254

**Discussion**

In this paper, we determined that age was the most critical indicator for predicting the severity of COVID-19. To improve the prediction accuracy, we proposed an MCDM algorithm, which combines the TOPSIS and NB classifier, to further select the indicators of patients' initial blood routine test. By ranking features, the MCDM algorithm selected 3 subsets including {Age, WBC, LYMC, NEUT}, {AGE WBC, LMYC} and {Age, NEUT, LYMC}, all of which can achieve 0.82 accuracy and 0.93 AUC.

Previous studies have shown that elderly COVID-19 patients with multiple concomitant diseases tend to develop Multiple Organ Failure (MOFE), which may lead to high morality in elderly patients infected by SARS-CoV-2. According to the latest meta-analysis of the elderly in the European community, the prevalence of frailty is around 15% for the elderly 65 years and older (*13*), and the case fatality rate of patients over 85 years old is 1,000 times that of patients aged 5-17 years (*14*). Our research indicated that age was the most important indicator for predicting the severity of COVID-19, with an accuracy 0.77 and an AUC of 0.92. However, some elderly patients had a good prognosis, so prognostic evaluation and medical decision-making based on age alone might not be accurate enough.

We found that WBC, LYMC and NEUT in initial blood routine test results other than age are also important for predicting the severity of COVID-19. Guo et al. (*15*) pointed out that the MuLBSTA score revealed that multi-lobar infiltrates, lymphocytes ≤0.8×109/L, bacterial infection, smoking status, hypertension, and age ≥60 years could help prognosticate outcomes in COVID-19 patients. The elevated WBC/NEUT is a basic sign of bacterial infection. Bacterial co-infection in COVID-19 patients may develop severe form of disease, complicating the clinical situation (*16-18*). The control and elimination of viruses depends on humoral immunity. Viral

278    infections usually lead to abnormal changes in the levels of lymphocyte subsets which further

279    impaired immune system functionality. The decrease of LYMC is the simplest and most intuitive

280    indicator to predict the humoral immune response, indicating that the patient's T cell function is

281    defective (*19-21*). The count of lymphocyte subsets (CD4+ and CD8+ T cell), especially CD8+

282    T cell, is directly proportional to the severity of COVID-19 (*22,23*).

283    Although logistic regression can determine the key indicator {Age}, and discrete analysis can

284    find a better subset {Age, WBC LYMC}, it is difficult to determine the best subset due to the

285    small sample size or multicollinearity. Previous studies used the MCDM algorithm to evaluate

286    diagnostic tests (*24*) and help doctors hasten COVID-19 treatment (*25*). As far as we know, this

287    is the first time the MCDM algorithm has been used to predict the severity of COVID-19. It first

288    uses TOPSIS for feature ranking, and then combines the NB classifier for feature selection. Even

289    if the sample size is small, the MCDM algorithm can select 3 effective subsets {Age, WBC,

290    LYMC}, {Age, WBC, LYMC, NEUT} and {Age, NEUT, LYMC}. The selection process is

291    visual and interpretable helping doctors find the features of the progress of emerging infectious

292    diseases early, to make faster and better prevention and treatment plans. We used the ROC curve

293    to evaluate the predictive value of the features selected by the MCDM algorithm. The results

294    showed that the MCDM algorithm can not only find all effective subsets, but also predict stably

295    and accurately.

296    Age (*26-29*), underlying diseases (*30*), systemic immune status (*31*), and blood test results can be

297    used as key features to predict the severity of COVID-19. Although these features can improve

298    the accuracy of prediction (84%~93%), the tests are time-consuming, expensive, and labor-

299    intensive. Our algorithm can select features from blood test results to achieve a prediction

300    accuracy of 82%. During the COVID-19 pandemic, it is more in line with clinical needs and is

301    easy to promote and use in areas with different medical levels.

302    The feature selection may have some limitations, because there were only 196 cases and all were

303    from China. In future, we would like to collect more data and conduct multi-center evaluations.

304    **Conclusion**

305    We defined feature selection as a MCDM problem so that the algorithm can provide a reference

306    for clinical practice. The concise features {Age, WBC/NEUT, LYMC} and high accuracy are

307    very conducive to rapid triage of COVID-19 patients. Using the most common blood routine test,

308    medical institutions could better determine the quarantine, hospital admission, ICU assignment

309    of COVID-19 patients. The MCDM algorithm can be used for small sample data sets, and the

310    prediction is accurate and stable, which will help establish a rapid response mechanism in the

311    early stage of emerging infectious disease outbreaks.

312

319    **Conflict of Interests:** The authors declare that they have no conflict of interests.

320    **Author Contributions:** Jiaqing Luo, Lingyun Zhou, and Shujin Guo designed the study. Shujin

321    Guo collected data. Jiaqing Luo, Bo Li and Yunyu Feng developed the algorithm. Jiaqing Luo,

322    Yunyu Feng and Lingyun Zhou edited the manuscript. Lingyun Zhou and Shujin Guo reviewed

323    the manuscript.

324    **Ethics:** The study was approved by the ethics committee of Sichuan Provincial People's

325    Hospital. Participant consent was not required.

326    **Data Availability Statement:** The data that support the findings of this study are available on

327    request from the corresponding author. The data are not publicly available due to privacy or

328    ethical restrictions.

329

330 **References**

331 1. World Health Organization. Coronavirus disease 2019 (COVID-19) situation report–51.

332     Geneva, Switzerland: World Health Organization; 2020. https://www.who.int/docs/default-

333     source/coronaviruse/situation-reports/20200311-sitrep-51-covid-19.pdf?sfvrsn=1ba62e57_10

334 2. Wang D, Hu B, Hu C, Zhu F, Liu X, Zhang J, et al. Clinical characteristics of 138

335     hospitalized patients with 2019 novel coronavirus–infected pneumonia in Wuhan,

336     China. JAMA. 2020;323:1061-9. doi:10.1001/jama.2020.1585

337 3. Yang X, Yu Y, Xu J, Shu H, Liu H, Wu Y, et al. Clinical course and outcomes of critically ill

338     patients with SARS-CoV-2 pneumonia in Wuhan, China: a single-centered, retrospective,

339     observational study. Lancet Respir Med. 2020;8:475-81. doi:10.1016/S2213-2600(20)30079-

340     5

341 4. Liang W, Liang H, Ou L, Chen B, Chen A, Li C, et al. Development and validation of a

342     clinical risk score to predict the occurrence of critical illness in hospitalized patients with

343     COVID-19. JAMA Intern Med. 2020; e202033. doi:10.1001/jamainternmed.2020.2033

344 5. Wu Z, McGoogan JM. Characteristics of and important lessons from the coronavirus disease

345     2019 (COVID-19) outbreak in China: summary of a report of 72 314 cases from the Chinese

346     Center for Disease Control and Prevention. JAMA. 2020; 323:1239-42.

347 6. Novel Coronavirus Pneumonia Emergency Response Epidemiology Team. The

348     epidemiological characteristics of an outbreak of 2019 novel coronavirus diseases (COVID-

349     19) in China [Chinese]. Zhonghua Liu Xing Bing Xue Za Zhi. 2020; 41:145–51.

350     doi:10.3760/cma.j.issn.0254-6450.2020.02.003

351 7. World Health Organization (2020). Clinical Management of Severe Acute Respiratory

352     Infection When Novel Coronavirus (nCoV) Infection is Suspected: Interim Guidance.

353     Available online at :https://www.who.int/docs/default-source/coronaviruse/clinical-

354     management-of-novel-cov. pdf  (accessed January20, 2020)

355  8.  Lin L, Li TS. Interpretation of "guidelines for the diagnosis and treatment of novel

356     coronavirus (2019-ncov) infection by the national health commission (trial version 5)".

357     Zhonghua Yi Xue Za Zhi. 2020;100: E001. doi:10.3760/cma.j.issn.0376-2491.2020.0001

358  9.  Chandrashekar G, Sahin F. A survey on feature selection methods. Comput Electr Eng. 2014;

359     40:16-28.

360  10. Yoon K. "A reconciliation among discrete compromise situations". J Oper Res Soc. 1987;

361     38:277-86. doi:10.1057/jors.1987.44.

362  11. Bunkley N. "Joseph Juran, 103, Pioneer in Quality Control, Dies". The New York Times.

363     2008.

364  12. McCallum A. "Graphical Models, Lecture2: Bayesian Network Represention". Retrieved 22

365     October, 2019.

366  13. Lippi G, Plebani M. Procalcitonin in patients with severe coronavirus disease 2019 (COVID-

367     19): a meta-analysis. CLIN CHIM ACTA. 2020;505:190.

368  14. Mo P, Xing Y, Xiao Y, Deng L, Zhao Q, Wang H, et al. Clinical characteristics of refractory

369     COVID-19 pneumonia in Wuhan, China. Clinical Infectious Diseases.2020.

370  15. Guo L, Wei D, WU Y, ZHOU M, ZHANG X, Li Q, et al. Clinical features predicting

371     mortality risk in patients with viral pneumonia: the MuLBSTA score. FRONT

372     MICROBIOL. 2019;10:2752.

373  16. Ma Y, Hou L, Yang X, Huang Z, Yang X, Zhao N, et al. The association between frailty and

374     severe disease among COVID-19 patients aged over 60 years in China: a prospective cohort

375     study. BMC MED. 2020; 18:1-8.

376  17. Stawicki SP, Jeanmonod R, Miller AC, Paladino L, Gaieski DF, Yaffee AQ, et al. The 2019–

377      2020 novel coronavirus (severe acute respiratory syndrome coronavirus 2) pandemic: A joint

378      american college of academic international medicine-world academic council of emergency

379      medicine multidisciplinary COVID-19 working group consensus paper. Journal of global

380      infectious diseases. 2020;12:47.

381  18. Lai CC, Shih TP, Ko WC, Tang HJ, Hsueh PR. Severe acute respiratory syndrome

382      coronavirus 2 (SARS-CoV-2) and corona virus disease-2019 (COVID-19): The epidemic and

383      the challenges. Int J Antimicrobial Agents. 2020;105924

384  19. Wang F, Nie J, Wang H, Zhao Q, Xiong Y, Deng L, .et al. Characteristics of peripheral

385      lymphocyte subset alteration in COVID-19 pneumonia. J INFECT DIS. 2020;221:1762-9.

386  20. Diao B, Wang C, Tan Y, Chen X, Liu Y, Ning L, et al. Reduction and functional exhaustion

387      of T cells in patients with coronavirus disease 2019 (COVID-19). FRONT

388      IMMUNO, 2020;11:827.

389  21. Mathew D, Giles JR, Baxter AE, Oldridge DA, Greenplate AR, Wu JE, .et al. Deep immune

390      profiling of COVID-19 patients reveals distinct immunotypes with therapeutic implication.

391      Science.2020;369(6508).

392  22. Pallotto C, Suardi LR, Esperti S, Tarquini R, Grifoni E, Meini S, et al. Increased CD4/CD8

393      ratio as a risk factor for critical illness in coronavirus disease 2019 (COVID-19): a

394      retrospective multicentre study. INFECT DIS-NOR. 2020; 52: 675-7.

395  23. Ganji A, Farahani I, Khansarinejad B, Ghazavi A, Mosayebi G. Increased expression of CD8

396      marker on T-cells in COVID-19 patients. BLOOD CELL MOL DIS. 2020;102437.

397  24. Sayan, M., Sarigul Yildirim, F., Sanlidag, T., Uzun, B., Uzun Ozsahin, D., & Ozsahin, I.

398      (2020). Capacity Evaluation of Diagnostic Tests For COVID-19 Using Multicriteria

399     Decision-Making Techniques. Computational and Mathematical Methods in Medicine, 2020.

400   25. Albahri OS, Al-Obaidi JR, Zaidan AA, Albahri AS, Zaidan BB, Salih MM, et al. Helping

401     doctors hasten COVID-19 treatment: Towards a rescue framework for the transfusion of best

402     convalescent plasma to the most critical patients based on biological requirements via ml and

403     novel MCDM methods. COMPUT METH PROG BIO. 2020;196:105617.

404   26. Guan WJ, Zhong NS. Clinical Characteristics of Covid-19 in China. Reply. N Engl J Med.

405     2020;382:1861-62. doi:10.1056/NEJMc2005203

406   27. Chen N, Zhou M, Dong X, Qu J, Gong F, Han Y, et al. Epidemiological and clinical

407     characteristics of 99 cases of 2019 novel coronavirus pneumonia in Wuhan, China: a

408     descriptive study. Lancet. 2020;395:507-13. doi:10.1016/S0140-6736(20)30211-7

409   28. Chen R, Liang W, Jiang M, Guan W, Zhan C, Wang T, et al. Risk Factors of Fatal Outcome

410     in Hospitalized Subjects With Coronavirus Disease 2019 From a Nationwide Analysis in

411     China. Chest. 2020;158:97-105. doi:10.1016/j.chest.2020.04.010

412   29. Zhou F, Yu T, Du R, Fan G, Liu Y, Liu Z, et al. Clinical course and risk factors for mortality

413     of adult inpatients with COVID-19 in Wuhan, China: a retrospective cohort study. Lancet.

414     2020;395:1054-1062. doi:10.1016/S0140-6736(20)30566-3

415   30. Wu C, Chen X, Cai Y, Zhou X, Xu S, Huang H, et al. Risk factors associated with acute

416     respiratory distress syndrome and death in patients with coronavirus disease 2019 pneumonia

417     in Wuhan, China. JAMA internal medicine. JAMA Intern Med. 2020;180:1-11.

418     doi:10.1001/jamainternmed.2020.0994

419   31. Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, et al. Clinical features of patients infected

420     with 2019 novel coronavirus in Wuhan, China. Lancet. 2020;395:497-506.

421     doi:10.1016/S0140-6736(20)30183-5