

1 Using Machine Learning of Clinical Data to Diagnose 2 COVID-19

3 Wei Tse Lit^{†,1}, Jiayan Ma^{†,1}, Neil Shende^{†,1}, Grant Castaneda¹, Jaideep Chakladar¹, Joseph C. Tsai¹,
4 Lauren Apostol¹, Christine O. Honda¹, Jingyue Xu¹, Lindsay M. Wong¹, Tianyi Zhang¹, Abby Lee¹,
5 Aditi Gnanasekar¹, Thomas K. Honda¹, Selena Z. Kuo², Michael Andrew Yu³, Eric Y. Chang⁴,
6 Mahadevan "Raj" Rajasekaran⁵, Weg M. Ongkeko^{*1}

7

8 ¹Department of Otolaryngology-Head and Neck Surgery, University of California San Diego, La Jolla, CA

9 92093; Research Service, VA San Diego Healthcare System San Diego, La Jolla, CA 92161, USA

10 ²Department of Medicine, Columbia University Medical Center, New York, NY 10032, USA

11 ³Department of Internal Medicine, Emory University School of Medicine, Atlanta, GA 30322, USA

12 ⁴Department of Radiology, University of California San Diego, La Jolla, CA 92093, USA; Radiology Service,
13 VA San Diego Healthcare System San Diego, La Jolla, CA 92161, USA

14 ⁵Department of Urology, University of California San Diego, La Jolla, CA 92093, US; Urology Service, VA San
15 Diego Healthcare System, San Diego, CA 92161, USA

16 * Correspondence: rongkeko@health.ucsd.edu

17 [†] Authors contributed equally

18 Received: date; Accepted: date; Published: date

19 **Abstract:** The recent pandemic of Coronavirus Disease 2019 (COVID-19) has placed severe stress
20 on healthcare systems worldwide, which is amplified by the critical shortage of COVID-19 tests. In
21 this study, we propose to generate a more accurate diagnosis model of COVID-19 based on patient
22 symptoms and routine test results by applying machine learning to reanalyzing COVID-19 data
23 from 151 published studies. We aimed to investigate correlations between clinical variables, cluster
24 COVID-19 patients into subtypes, and generate a computational classification model for
25 discriminating between COVID -19 patients and influenza patients based on clinical variables
26 alone. We discovered several novel associations between clinical variables, including correlations

27 between being male and having higher levels of serum lymphocytes and neutrophils. We found
28 that COVID-19 patients could be clustered into subtypes based on serum levels of immune cells,
29 gender, and reported symptoms. Finally, we trained an XGBoost model to achieve a sensitivity of
30 92.5% and a specificity of 97.9% in discriminating COVID-19 patients from influenza patients. We
31 demonstrated that computational methods trained on large clinical datasets could yield ever more
32 accurate COVID-19 diagnostic models to mitigate the impact of lack of testing. We also presented
33 previously unknown COVID-19 clinical variable correlations and clinical subgroups.

34 **Keywords:** COVID-19, machine learning, diagnostic model

35

36 **1. Introduction**

37 COVID-19 is a severe respiratory illness caused by the virus SARS-CoV-2. The scientific
38 community has focused on this disease with near unprecedented intensity. However, the majority of
39 primary studies published on COVID-19 suffered from small sample sizes[1, 2]. While a few primary
40 research studies reported on dozens or hundreds of cases, many more studies reported on less than
41 20 patients[3, 4]. Therefore, there is an urgent need to collate all available published data on the
42 clinical characteristics of COVID-19 from different studies to construct a comprehensive dataset for
43 gaining insights into the pathogenesis and clinical characteristics of COVID-19. In this study, we aim
44 to perform a large-scale meta-analysis to synthesize all published studies with COVID-19 patient
45 clinical data, with the goal of uncovering novel correlations between clinical variables in COVID-19
46 patients. We will then apply machine learning to reanalyze the data and construct a computational
47 model for predicting whether someone has COVID-19 based on their clinical information alone.

48 We believe that the ability of predicting COVID-19 patients based on clinical variables and
49 using an easily accessible computational model would be extremely useful to address the
50 widespread lack of testing capabilities for COVID-19 worldwide. Because many countries and
51 hospitals are not able to allocate sufficient testing resources, healthcare systems are deprived of one

52 of their most effective tools for containing a pandemic: identification of case hotspots and targeted
53 action towards regions and specific individuals with the disease[5]. The scale of the testing shortage
54 calls for methods for diagnosing COVID-19 that use resources local healthcare facilities currently
55 have. We propose the development of a disease prediction model based on clinical variables and
56 standard clinical laboratory tests.

57 A number of meta-analyses have been done on COVID-19, but almost none of them
58 comprehensively included data from all published studies. Three different meta-analyses, published
59 in February, March, and April of 2020, included data from 10, 8, and 31 articles, respectively[6-8]. We
60 included 151 articles, comprising 413 patients, in our analysis. To the best of our knowledge, no
61 study has performed a large-scale machine learning analysis on clinical variables to obtain a
62 diagnostic model. We believe that our study will be an important step towards leveraging the full
63 extent of published clinical information on COVID-19 patients to inform diagnosis of COVID-19,
64 instead of relying on general guidelines for symptoms that do not take into account the association
65 between different clinical variables.

66 2. Results

67 2.1 Compilation of patient data and summary of clinical variables

68 After compiling information from 151 published studies, we present a total of 42 different
69 clinical variables, including 21 categorical and 21 continuous variables, that are reported in more
70 than 1 study. Discrete variables include nominal categorical variables like gender, which is 49.49%
71 (194 patients) male and 50.51% (198 patients) female, and ordinal categorical variables like
72 lymphocytes level, of which 86 patients (48.86%) have low levels, 73 patients (41.47%) have normal
73 levels, and 17 patients (9.65%) have high levels. Continuous variables include age, which has a mean
74 of 38.91 years and variance 21.86 years, and serum neutrophil levels, which has a mean of 6.85×10^9
75 cells/L and a variance of 12.63×10^9 cells/L. Certain variables, including all counts of blood cell
76 populations, have both ordinal and continuous components. The continuous component describes
77 the raw count of these populations, while the ordinal component describes whether these counts are

78 within normal range, below normal range, or above normal range. Summary for all data are shown
79 in Table 1.

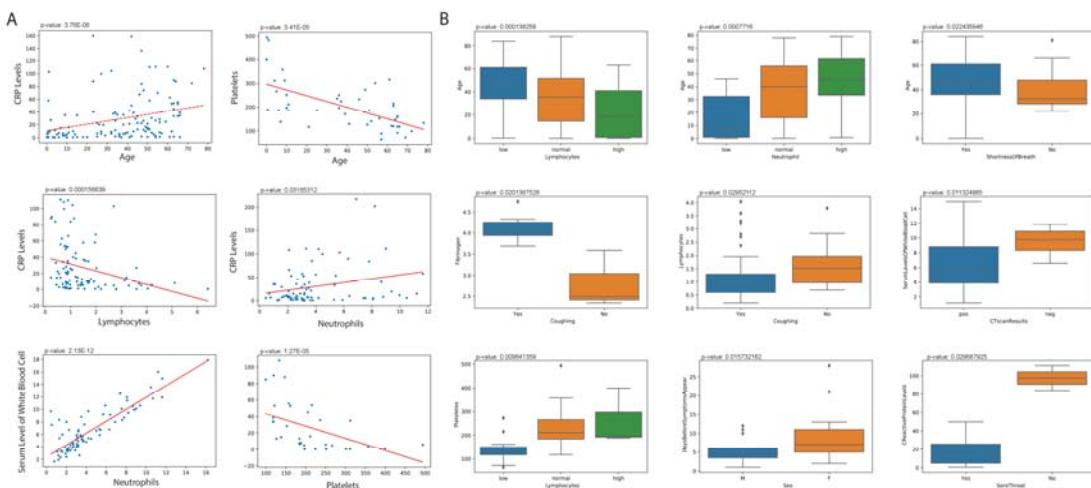
80

81 *2.2 Relationship between pairs of clinical variables*

82 We performed correlation between all possible pairs of clinical variables to uncover potentially
83 important associations (Table S1). If both variables are continuous, the Spearman correlation test is
84 applied ($p<0.05$). Among 143 Spearman correlation tests, 27 show significant correlation, with 9 that
85 involve age, corroborating reports that age plays a critical role in the development of COVID-19[9].

86 We observed C reactive protein (CRP) levels and serum platelets levels to be the variables with the
87 strongest correlations with age (Figure 1A). CRP levels, an indicator of inflammation, are positively
88 correlated with increasing age, while platelets levels are negatively correlated with increasing age.
89 Other than age, we observed a negative correlation between the levels of CRP and lymphocytes
90 levels and a positive correlation between CRP levels and neutrophil levels (Figure 1A). This result
91 suggests that inflammation is most likely driven by neutrophils. The serum levels of white blood
92 cells are also strongly correlated with neutrophils, further suggesting that white blood cell counts
93 are heavily influenced by neutrophil levels (Figure 1A).

94

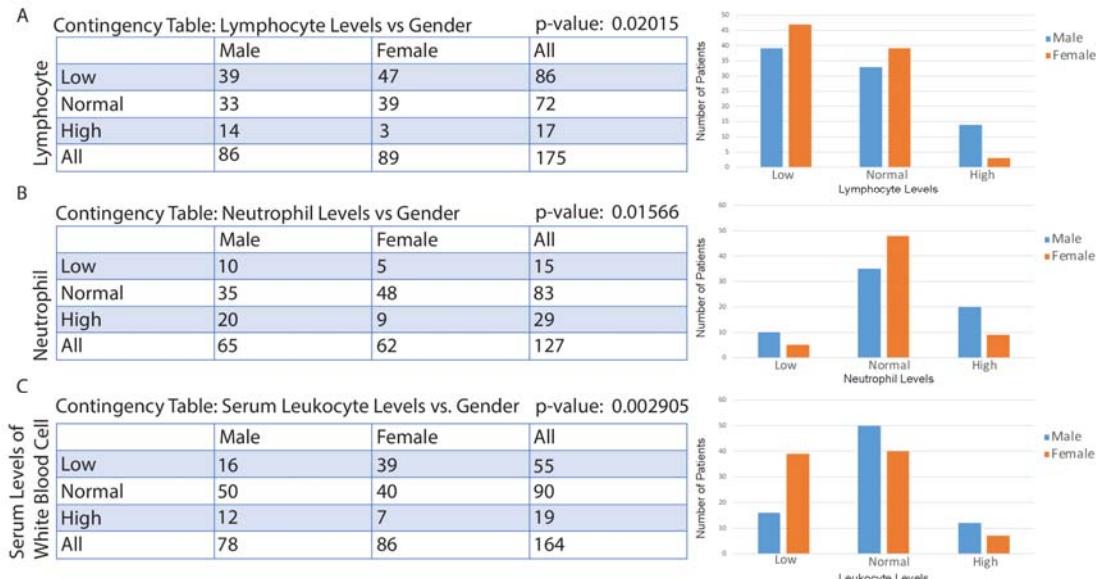


95 **Figure 1.** Select correlations with continuous clinical variables for COVID-19 patients. **A.**
96 Correlations between two continuous variables (Spearman, $p<0.05$). **B.** Correlations between one
97 continuous and one categorical variable (Kruskal-Wallis test, $p<0.05$).

98

99 For pairs of one continuous and one discrete variable, the Kruskal-Wallis test is applied
100 ($p<0.05$). Among 319 Kruskal-Wallis comparisons, 36 correlations were significant. Some of the
101 significant pairs overlapped with correlations between two continuous variables for variables that
102 have both ordinal and continuous components. Such correlations are not displayed twice in Figure 1.
103 We found again that age correlated significantly with multiple variables, including negative
104 correlation with lymphocyte levels, positive correlation with neutrophil levels, and positive
105 correlation with shortness of breath (Figure 1B). Other interesting associations were also discovered.
106 Coughing was found to be correlated with increasing fibrinogen levels and decreasing lymphocyte
107 levels. Those with lower levels of serum white blood cells (leukocytes) are more likely to report a
108 positive CT scan result for pneumonia. Females may experience a greater number of days before
109 symptoms appear. Finally, we found that sore throat decreases with increasing CRP levels (Figure
110 1B).

111 For pairs of two categorical variables, a two-tail chi-square test is applied. 42 out of 309
112 comparisons showed significant correlation, with few overlaps with former tests. Gender is involved
113 in 6 of the significant correlation, indicating significant gender differences in COVID-19.
114 Contingency tables of selected significant correlations are shown in Figure 2. Males were found to
115 have higher lymphocyte and neutrophil levels than females (Figure 2A,B). Females were found to be
116 more likely to have lower levels of serum white blood cells (Figure 2C).



117

118 **Figure 2.** Correlations between gender and another categorical variable. A. Correlation between
119 lymphocyte level categories and gender. B. Correlation between neutrophil level categories and
120 gender. C. Correlation between serum leukocyte level categories and gender. A contingency table
121 and a bar plot of the number of patients in each level are displayed for each correlation.

122

123 2.3 Clustering of patients into subcategories of COVID-19

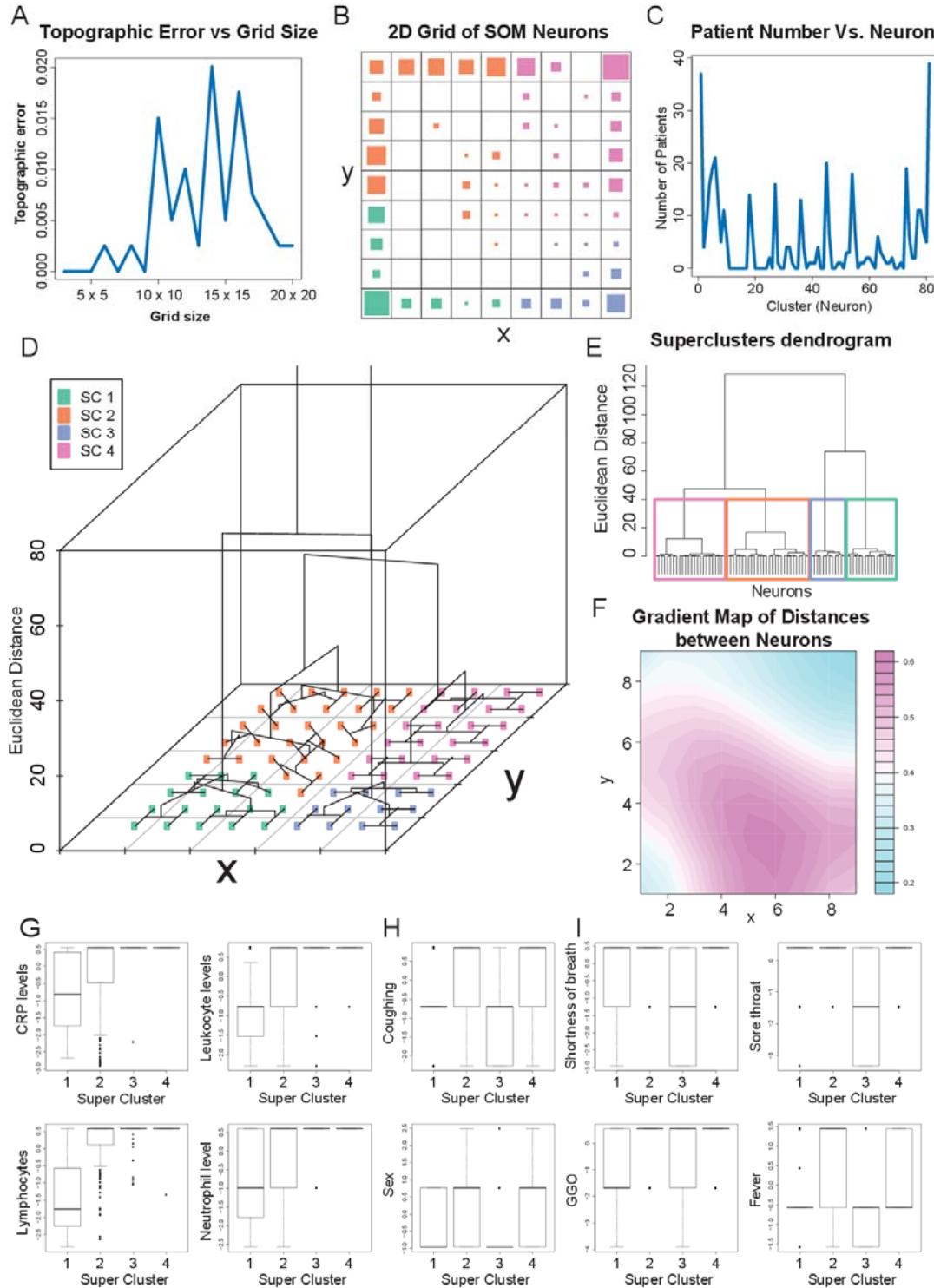
124 We next aim to cluster COVID-19 patients based on clinical variables using machine learning.

125 We chose the well-known SOM algorithm for clustering. SOM is a neural network that has a set of
126 neurons organized on a 2D grid[10]. All neurons are connected to all input units (individual
127 patients) by a weight vector. The weights are determined through iterative evaluations of a Gaussian
128 neighborhood function, with the result of creating a 2D topology of neurons to model the similarity
129 of input units (individual patients). The algorithm outputs a map that assigns each sample to one of
130 the neurons on the 2D grid, with samples in the same neuron being the most similar to one another.

131 Similarity of samples decrease with distance between neurons on the 2D map. Missing variables will
132 be ignored from the SOM model when deriving a neural topology.

133 We generated square SOM neuron grids with side lengths 3 through 20 using the trainSOM
134 function in the R package SOMbrero. The grids with side lengths 3, 4, 5, 7, and 9 all had topographic

135 errors of 0 (Figure 3A). Of these, we chose the biggest grid (9x9=81 clusters) as our model. After the
136 patients were assigned to neurons, an analysis of variance (ANOVA) test was performed to test
137 which variables actively participate in the clustering. Of the 48 clinical variables we inputted, 27
138 were found to have very high significance ($p<0.001$) (Table S2). We then reran the SOM using the
139 27 variables on a 9x9 grid. This grid is displayed on Figure 3B and has a final energy of 8.139248. The
140 largest neuron has 39 patients, the second largest has 37, the third largest has 21, and the fourth
141 largest has 20 (Figure 3C).



142

143 **Figure 3.** Summary of COVID-19 patient clustering using SOM. **A.** Plot of topographic error of the
144 2D SOM grid vs. size of the grid. **B.** 2D plot of SOM neurons after retaining only the most
145 significant clinical variable for analysis. Each small grid represents a neuron, and the size of the

146 square in each grid represents the number of patients associated with each neuron. The color code
147 corresponds to superclusters presented in panel D. C. Plot of number of patients in each neuron.
148 D. 3D dendrogram summarizing the neurons into superclusters. E. 2D dendrogram with the same
149 information as the dendrogram in part D. In both dendograms, the vertical axis represents the
150 relative distance between clusters, which can be known between any two clusters by looking at
151 the branch point where they diverge. F. Gradient map where light blue regions of the SOM depict
152 higher similarity of neurons with each other. G. Boxplots of immune-associated clinical variables
153 that differentiate superclusters. H. Boxplots in which superclusters 1 and 3 display similar trends.
154 I. Boxplots in which only one supercluster has a median at a different value from the other three.
155 All variables have been previously normalized. For binary variables, only three possible positions
156 on the vertical axis is possible: the bottom one being no, the middle one being yes, and the top one
157 being missing. For the gender (sex) variable, the bottom position is female, the middle is male,
158 and the top one is missing.

159

160 *2.4 Clinical characteristics of COVID-19 clusters*

161 We then examined the defining features of patients assigned to the same neurons. We
162 investigated four neurons associated with the largest number of patients and identified the four
163 variables with the smallest nonzero standard deviations for each patient cluster. In the largest
164 cluster, with 39 patients, the four smallest nonzero standard deviations were for the variables region
165 of infection, sore throat, RT-PCR results, and coughing. In the second largest cluster, with 37
166 patients, the variables were baby death if pregnant, lymphocyte levels, fever, and coughing. In the
167 next largest cluster, with 21 patients, the variables were sore throat, duration of illness in days,
168 RT-PCR results, and coughing. In the fourth largest cluster, with 20 patients, the variables were sore
169 throat, fever, coughing, and age.

170 We next used the function superClass to compute the relative Euclidean distances between the
171 81 patient clusters and form superclusters. The relative distances between the individual clusters are
172 shown in Figure 3D-E. We divided the 81 clusters into 4 superclusters, which are represented in
173 Figure 3B by the color of the squares. Supercluster 1 was formed with 24 neurons, supercluster 2 had

174 28 neurons, supercluster 3 had 12, and supercluster 4 had 17 neurons. Visualizing distances between
175 neighboring neurons, we found that the distances are the smallest at corners of the grid, especially
176 the upper right-hand corner (Figure 3F). This corner corresponds to supercluster 4, suggesting that
177 patients within this cluster may be especially similar.

178

179 Next, we sought to determine the clinical features that effectively distinguish these
180 superclusters. We performed Kruskal-Wallis testing on the values of the 27 variables across the four
181 superclusters. 24 variables were significantly different between the superclusters ($p < 0.05$) (Table S3).
182 We discovered that the clinical variables exhibit 3 main types of correlations with the superclusters:
183 continuous increase in value from cluster 1 to cluster 4 (Figure 3G), clusters 1 and 3 exhibiting the
184 same distribution and clusters 2 and 4 exhibiting another distribution (Figure 3H), and 3 clusters
185 exhibiting the same median (Figure 3I). From these analyses, we could infer that patients with low
186 levels of CRP and serum immune cells likely define cluster 1. Cluster 1 patients are also
187 predominantly female. Cluster 2 contains patients with slightly higher levels of CRP and serum
188 immune cells than cluster 1. Compared to cluster 1 patients, fewer cluster 2 patient report coughing
189 and fever. Cluster 2 patients are predominantly male. Cluster 3 contains patients with few reported
190 symptoms, including less coughing, shortness of breath, fever, and sore throat. Cluster 3 is
191 overwhelmingly female. Cluster 4 most likely contains patients not belonging to the other 3 clusters
192 as it has few distinguishing features and high levels of missing data.

193

194 *2.5 Creation of a diagnostic model for COVID-19 based on clinical variables*

195 Because it can be difficult to distinguish influenza from COVID-19, we downloaded clinical
196 data collected for influenza from a study by Cheng et. al. and from the Influenza Research
197 Database[11, 12]. Machine learning was then used to perform a classification task to discriminate
198 between influenza and COVID-19. For machine learning, we employed the algorithm Extreme
199 Gradient Boosting (XGBoost) using Python. XGBoost is a novel, state-of-the-art machine learning

200 algorithm that has been shown to outperform other more traditional algorithms in its accuracy and
201 efficiency[13]. It can also take both continuous and discrete inputs and handle sparse data, in
202 addition to having highly optimizable hyper-parameters[14].

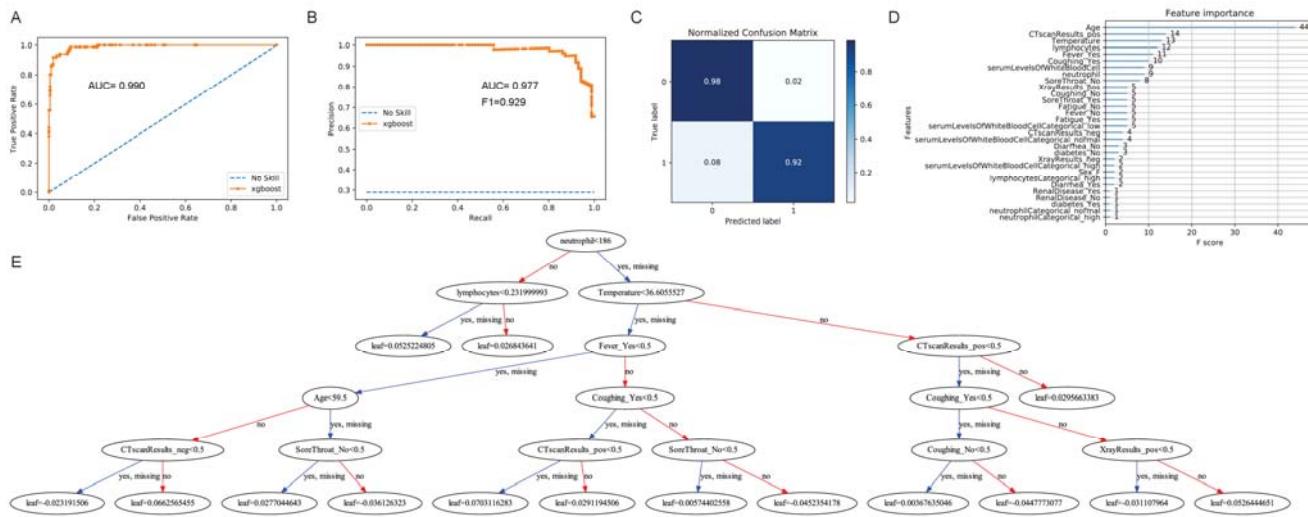
203 The datasets from non-COVID patients and COVID-19 patients were merged and then split into
204 training and testing patient sets, with 80% and 20% of the patients, respectively. Categorical
205 variables were encoded as dummy variables. We then tuned the model using the Bayesian
206 optimization method for hyperparameter search. We found the best hyperparameters to be
207 gamma=0.0933, learning rate= 0.4068, max depth=6.558, and n_estimators=107.242.

208

209 *2.6 Evaluation of XGBoost classification outcomes*

210 From the ROC curve of prediction results, we obtained an AUC of 0.990 (Figure 4A). However,
211 because there is an imbalance of class in our input (i.e. we have significantly more influenza patients
212 than COVID-19 patients), the precision recall (PR) curve may be better able to present our model's
213 results. ROC curves could be significantly influenced by skewing the distribution of classes in
214 classification, while PR curves would not be impacted by this action. We observed a slightly lower
215 AUC of 0.977 in our PR curve and computed the F1 score to be 0.929 (Figure 4B), which suggest that
216 our model is still highly accurate even when class imbalances are taken into account. The prediction
217 result from XGBoost's predict function was used to plot a confusion matrix (Figure 4C). From the
218 confusion matrix, we calculated a sensitivity of 92.5% and a specificity of 97.9%. We found the most
219 important features in our prediction model to be age, CT scan result, temperature, lymphocyte
220 levels, fever, and coughing, in order of decreasing importance (Figure 4D). We also provided a
221 6-level decision tree sample of our XGBoost model (Figure 4E), which is not a representation of our
222 full model.

223



224
225
226
227

Figure 4. Summary of XGBoost classification of COVID-19 and influenza patients. **A.** ROC curve of prediction. **B.** Precision recall curve of prediction. **C.** Confusion matrix of prediction. **D.** Variables most important for classification, listed by decreasing order of importance. **E.** 6-level sample model of SOM decision tree construction.

228 3. Discussion

229 As the recent pandemic of COVID-19 unfolds across the world, the inability of countries to test
230 their citizens is heavily impacting their healthcare system's ability to fight the epidemic. Testing is
231 necessary for the identification and quarantine of COVID-19 patients. However, the multi-step
232 process required for the conventional SARS-CoV-2 test, via quantitative polymerase chain reaction
233 (qPCR), is creating difficulties for countries to test large numbers of suspected patients[15]. Testing
234 begins with a healthcare worker taking a swab from the patient. The swab is sent to a laboratory, and
235 viral RNA is extracted from the sample and reverse transcribed into DNA. The DNA is tagged with
236 a fluorescent dye and then amplified using a qPCR machine. If a high level of fluorescence is
237 observed compared to control, the sample is positive with SARS-CoV-2. Each step of the testing
238 process is susceptible to severe shortages[16].

239 In this study, we aim to mine published clinical data of COVID-19 patients to generate a new
240 diagnostic framework. We hypothesize that novel or complex associations between clinical variables

241 could be exploited for diagnosis with the aid of machine learning. Not only may underlying
242 relationships between clinical variables in COVID-19 be useful for the development of a
243 computational diagnostic test based on signs, symptoms, and laboratory results, these correlations
244 can also yield critical insights into the biological mechanisms of COVID-19 transmission and
245 infection.

246 Using correlational tests, we corroborated previous findings and expected results for
247 COVID-19 patients but also uncovered novel relationships between clinical variables. We found that
248 age is correlated with CRP level, an indicator of inflammation, and decreased platelet levels. It is
249 known that as age increases, the proinflammatory response becomes stronger, leading to increasing
250 CRP and decreasing platelet levels[17]. However, we found surprising correlations with gender,
251 including higher serum neutrophil and leukocyte levels in males compared to females. According to
252 the National Health and Nutrition Examination Survey, with data from over 5,600 individual, few
253 differences exist between male and females in the serum levels of these cells[18]. Another study with
254 200 samples found that neutrophils are generally higher in women[19]. Correlations with gender
255 observed here may offer a piece of the explanation for why men infected with COVID-19 seem to
256 experience poorer prognosis, one of the important outstanding questions of COVID-19[20].

257 We also classified COVID-19 patients into different clusters using the SOM machine learning
258 algorithm. Two of the clusters are defined by low vs. high levels of immunological parameters,
259 including immune cell counts and CRP levels. A third cluster is defined by a tendency for less
260 reported symptoms, including sore throat, fever, and shortness of breath, and is predominantly
261 female.

262 Finally, using the machine learning algorithm XGBoost, we constructed a computational model
263 that successfully classified influenza patients from COVID-19 patients with high sensitivity and
264 specificity. We believe that our model demonstrated the feasibility of using data mining and
265 machine learning to inform diagnostic decisions for COVID-19. Such a model could be extremely

266 useful for more effective identification of COVID-19 cases and hotspots, which could allow health
267 officials to act before testing shortages could be addressed.

268 Despite promising results, several limitations exist for our study, all of which stem from the
269 lack of large-scale clinical data. First, our sample size is severely limited because most clinical reports
270 published do not publish individual-level patient data. Second, data on influenza signs and
271 symptoms are equally inaccessible. We were only able to locate data for patients with H1N1
272 influenza A, which is not one of the active strains in the current influenza season. Third, many of our
273 data sources are case studies that focused on specific cohorts of COVID-19 patients. This increases
274 the chance of us capturing a patient population that is not representative of the general population,
275 although this is an inherent risk of sampling. We anticipate that as more data are made openly
276 available in the weeks and months to come, we will be able to build a more robust computational
277 model. Therefore, we intend to provide the model we constructed as a computational framework for
278 computation-aided diagnosis of COVID-19 data rather than a ready-to-use model. We also
279 encourage researchers around the world to release de-identified patient data to aid in data mining
280 and machine learning efforts against COVID-19.

281

282 **4. Materials and Methods**

283 *4.1 Literature search and inclusion criteria for studies*

284 Patient clinical data were manually curated from a PubMed search with the keyword
285 "COVID-19." A total of 1,439 publications, dating from January 17, 2020 to March 23, 2020, were
286 reviewed. All publications with no primary clinical data, including reviews, meta-analyses, and
287 editorials, were excluded from our analysis. After manual review, we found 151 studies with
288 individual-level data, encompassing data from 413 patients. All individual patient data with 2 or more
289 clinical variables reported per patient were included. Clinical variables sought for included
290 demographics, signs and symptoms, laboratory test results, imaging results, and COVID-19

291 diagnosis. No formal review protocol was used. No bias was assessed within or across studies
292 because we did not include any clinical trials or case-control studies.

293 For our machine learning classification task to discriminate COVID-19 patients from influenza
294 patients, we used clinical variables for 21 influenza patients from a study by Cheng et. al. and 1050
295 patients from the Influenza Research Database[11, 12]. Only H1N1 Influenza A virus cases were
296 included because of difficulties locating data from other strains.

297

298 *4.2 Correlational tests between pairs of clinical variables*

299 We sought to uncover correlations that could yield critical insights into the clinical
300 characteristics of COVID-19 by correlating every variable to each other. For two continuous
301 variables, the Spearman correlation test was applied. For one continuous variable and one
302 categorical variable, the Kruskal-Wallis test was applied. For two different categorical variables, the
303 chi-squared test will be applied. All statistical tests were considered significant if the *p*-value is 0.05
304 or below.

305

306 *4.3 Machine learning for classification of COVID-19 patients into subtypes*

307 A self-organizing map (SOM) is an artificial neural network that constructs a two-dimensional,
308 discretized depiction (map) of the training set. We used the SOM algorithm to cluster our patients
309 based on similar patterns of clinical variables. The SOMbrero package in R was used[21]. Because
310 clustering of neurons are performed using Euclidean Distance, we first standardized each clinical
311 variable to ensure that they are equally weighted.

312 The trainSOM function was used to implement numeric SOM on our data set, which is inputted
313 as an N x P matrix, with N=398 patients and P=48 variables. From this, we selected 27 clinical
314 variables with very high significance (*p*<0.001) after running an ANOVA test across all neurons
315 and ran another iteration of trainSOM with these variables. We generated SOMs from the 3x3
316 neuron grid to 20x20 neuron grid and selected the 9x9 SOM with 81 neurons as our final model

317 based on minimal topographic error. We then aggregated the neurons into super-clusters using
318 superClass method in SOMbrero.

319

320 *4.4 Preprocessing of data for machine learning classification*

321 Data were preprocessed by combining data from COVID-19 cases and influenza cases into a
322 single matrix, followed by removal of any clinical variables that were not present in both the
323 COVID-19 dataset and the influenza dataset. 19 clinical variables were included as machine learning
324 input. The variables include age, sex, serum levels of neutrophil (continuous and ordinal), serum
325 levels of leukocytes (continuous and ordinal), serum levels of lymphocytes (continuous and ordinal),
326 result of CT scans, result of chest X-rays, reported symptoms (diarrhea, fever, coughing, sore throat,
327 nausea, and fatigue), body temperature, and underlying risk factors (renal diseases and diabetes).

328 Categorical data were converted to dummy variables using the get dummies function in Pandas
329 because non-numerical data are not allowed in our machine learning algorithm.

330 *4.5 Performing XGBoost classification*

331 The eXtreme Gradient Boosting algorithm (XGBoost), an ensemble machine learning method
332 widely known for its superior performance over other machine learning methods, was selected for
333 our study[13]. We first split our data into 80% training dataset and 20% testing dataset. 5-fold
334 cross-validation was then performed, with 70 boosting rounds (iterations), and fed into a Bayesian
335 optimization function for calculation of the best hyperparameters for XGBoost. The
336 hyperparameters tuned included max depth, gamma, learning rate, and n estimators. Bayesian
337 optimization was performed with an initial 8 steps of random exploration followed by 5 iterations.
338 The expected improvement acquisition function was used.

339

340 *4.6 Evaluation of XGBoost results*

341 XGBoost results were evaluated by plotting a receiver operating characteristic (ROC) curve and
342 a precision recall (PR) curve. The area under the curve (AUC) was also calculated for both curves.

343 **Declarations**

344 *Ethics approval and consent to participate:* Not applicable.

345 *Consent for publication:* Not applicable.

346 *Availability of data and materials:* The datasets during and/or analysed during the current study
347 available from the corresponding author on reasonable request.

348 *Competing Interests:* The authors declare that they have no competing interests.

349 *Funding:* University of California, Office of the President/Tobacco-Related Disease Research
350 Program Emergency COVID-19 Research Seed Funding Grant (R00RG2369) to W.M.O.

351 *Author Contributions:* Conceptualization, W.M.O.; methodology, W.M.O. and W.T.L.; software,
352 W.T.L. J.M., and N.S.; validation, W.T.L., J.M., and N.S.; formal analysis, W.T.L., J.M., and N.S.;
353 investigation, W.T.L., J.M., N.S., G.C., J.C., J.C.T., L.A., C.O.H., J.X., L.M.W., T.Z., A.L., A.G., and
354 T.K.H.; resources, W.M.O.; data curation, W.T.L., J.M., N.S., G.C., J.C., J.C.T., L.A., C.O.H., J.X.,
355 L.M.W., T.Z., A.L., A.G., and T.K.H.; writing—original draft preparation, W.T.L. and N.S.;
356 writing—review and editing, W.M.O., E.Y.C., M.R.R., S.Z.K., and M.A.Y.; visualization, W.T.L., J.M.,
357 and N.S.; supervision, W.M.O.; project administration, W.M.O.; funding acquisition, E.Y.C. All
358 authors have read and agreed to the published version of the manuscript.

359 *Acknowledgments:* Not applicable.

360 **Abbreviations**

CRP C-reactive Protein

ANOVA Analysis of Variance

SOM Self-organizing map

XGBoost Extreme Gradient Boosting

ROC Receiver Operating Characteristic

AUC Area Under the Curve

PR Precision-Recall

361 References

- 362 1. Chang Mo G, Yuan X, Tao Y, Peng X, Wang F, Xie L, Sharma L, Dela Cruz CS, Qin E: **Time Kinetics of**
363 **Viral Clearance and Resolution of Symptoms in Novel Coronavirus Infection.** *Am J Respir Crit Care*
364 *Med* 2020.
- 365 2. Zhang MQ, Wang XH, Chen YL, Zhao KL, Cai YQ, An CL, Lin MG, Mu XD: **[Clinical features of 2019**
366 **novel coronavirus pneumonia in the early stage from a fever clinic in Beijing].** *Zhonghua Jie He He*
367 *Hu Xi Za Zhi* 2020, **43**(3):215-218.
- 368 3. Feng K, Yun YX, Wang XF, Yang GD, Zheng YJ, Lin CM, Wang LF: **[Analysis of CT features of 15**
369 **Children with 2019 novel coronavirus infection].** *Zhonghua Er Ke Za Zhi* 2020, **58**(0):E007.
- 370 4. Li Y, Guo F, Cao Y, Li L, Guo Y: **Insight into COVID-2019 for pediatricians.** *Pediatr Pulmonol* 2020.
- 371 5. HUANG P: **If Most Of Your Coronavirus Tests Come Back Positive, You're Not Testing Enough.** In:
372 *NPR*. 2020.
- 373 6. Sun P, Qie S, Liu Z, Ren J, Li K, Xi J: **Clinical characteristics of hospitalized patients with**
374 **SARS-CoV-2 infection: A single arm meta-analysis.** *J Med Virol* 2020.
- 375 7. Yang J, Zheng Y, Gou X, Pu K, Chen Z, Guo Q, Ji R, Wang H, Wang Y, Zhou Y: **Prevalence of**
376 **comorbidities in the novel Wuhan coronavirus (COVID-19) infection: a systematic review and**
377 **meta-analysis.** *Int J Infect Dis* 2020.
- 378 8. Cao Y, Liu X, Xiong L, Cai K: **Imaging and Clinical Features of Patients With 2019 Novel**
379 **Coronavirus SARS-CoV-2: A systematic review and meta-analysis.** *J Med Virol* 2020.
- 380 9. Kolifarhood G, Aghaali M, Mozafar Saadati H, Taherpour N, Rahimi S, Izadi N, Hashemi Nazari SS:
381 **Epidemiological and Clinical Aspects of COVID-19; a Narrative Review.** *Arch Acad Emerg Med* 2020,
382 **8**(1):e41.
- 383 10. Jerez JM, Molina I, Garcia-Laencina PJ, Alba E, Ribelles N, Martin M, Franco L: **Missing data**
384 **imputation using statistical and machine learning methods in a real breast cancer problem.** *Artif*
385 *Intell Med* 2010, **50**(2):105-115.
- 386 11. Cheng Y, Zhao H, Song P, Zhang Z, Chen J, Zhou YH: **Dynamic changes of lymphocyte counts in**
387 **adult patients with severe pandemic H1N1 influenza A.** *J Infect Public Health* 2019, **12**(6):878-883.
- 388 12. Squires RB, Noronha J, Hunt V, Garcia-Sastre A, Macken C, Baumgarth N, Suarez D, Pickett BE, Zhang
389 Y, Larsen CN *et al:* **Influenza research database: an integrated bioinformatics resource for influenza**
390 **research and surveillance.** *Influenza Other Respir Viruses* 2012, **6**(6):404-416.
- 391 13. Chen T, Carlos G: **XGBoost: A Scalable Tree Boosting System.** *KDD '16: Proceedings of the 22nd ACM*
392 *SIGKDD International Conference on Knowledge Discovery and Data Mining* 2016:9.
- 393 14. Al'Aref SJ, Maliakal G, Singh G, van Rosendael AR, Ma X, Xu Z, Alawamli OAH, Lee B, Pandey M,
394 Achenbach S *et al:* **Machine learning of clinical variables and coronary artery calcium scoring for the**
395 **prediction of obstructive coronary artery disease on coronary computed tomography angiography:**
396 **analysis from the CONFIRM registry.** *Eur Heart J* 2020, **41**(3):359-367.
- 397 15. Hollingsworth J: **A coronavirus test can be developed in 24 hours. So why are some countries still**
398 **struggling to diagnose?** In: *CNN*. 2020.
- 399 16. Yong E: **How the Pandemic Will End.** In: *The Atlantic*. 2020.
- 400 17. Molloy EJ, Bearer CF: **COVID-19 in children and altered inflammatory responses.** *Pediatr Res* 2020.
- 401 18. Andersen CJ, Vance TM: **Gender Dictates the Relationship between Serum Lipids and Leukocyte**
402 **Counts in the National Health and Nutrition Examination Survey 1999(-)2004.** *J Clin Med* 2019, **8**(3).

- 403 19. Bain BJ, England JM: **Normal haematological values: sex difference in neutrophil count.** *Br Med J*
404 1975, **1**(5953):306-309.
405 20. Wenham C, Smith J, Morgan R, Gender, Group C-W: **COVID-19: the gendered impacts of the**
406 **outbreak.** *Lancet* 2020, **395**(10227):846-848.
407 21. Boelaert J, Bendhaiba L, Olteanu M, Villa-Vialaneix N: **SOMbrero: an R Package for Numeric and**
408 **Non-numeric Self-Organizing Map.** 2013.

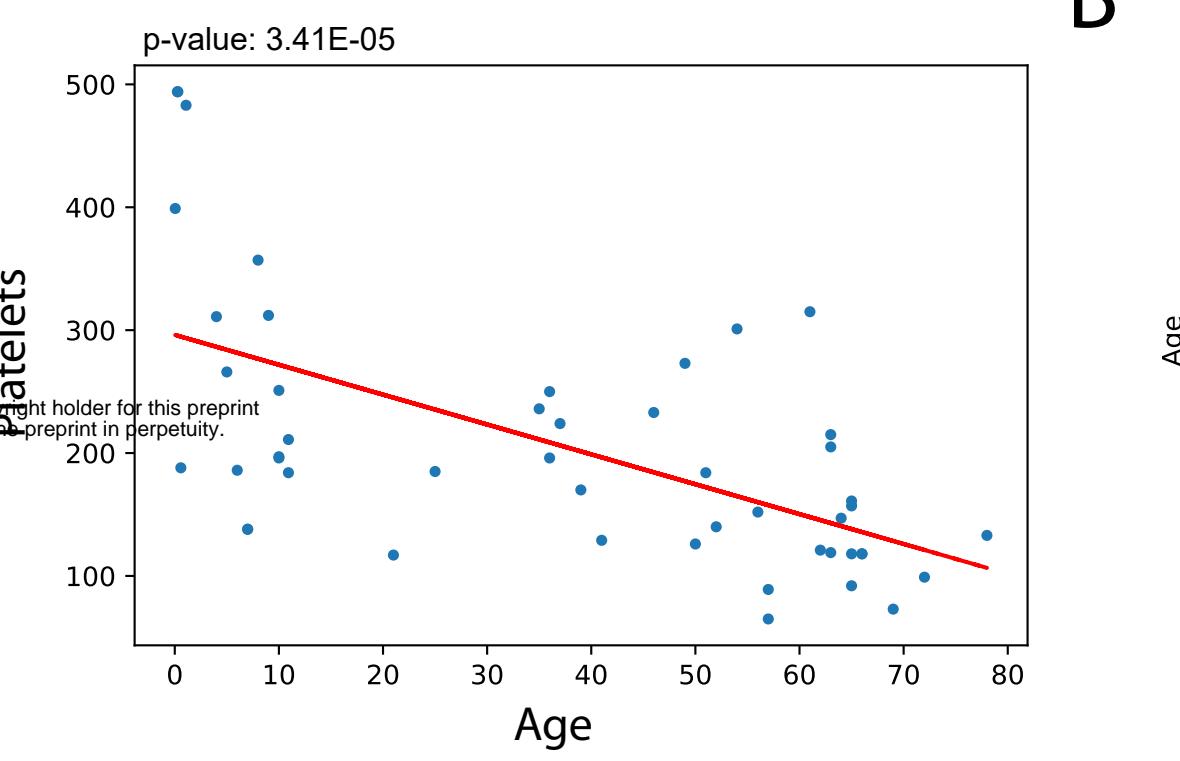
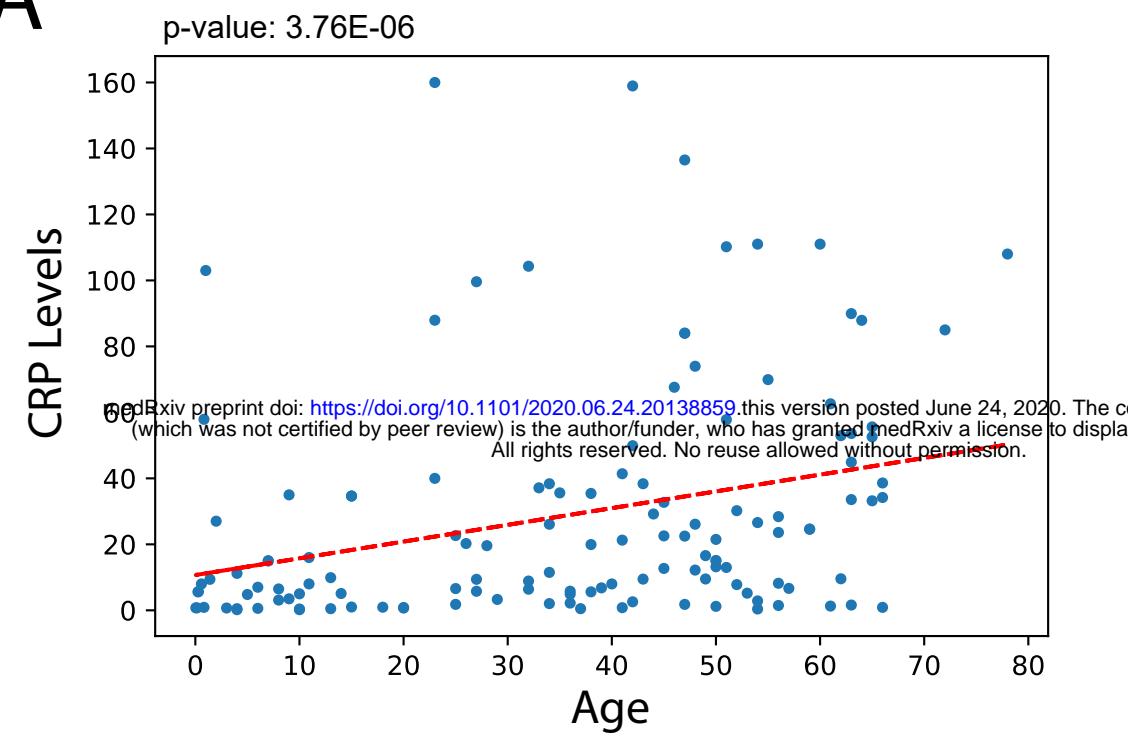
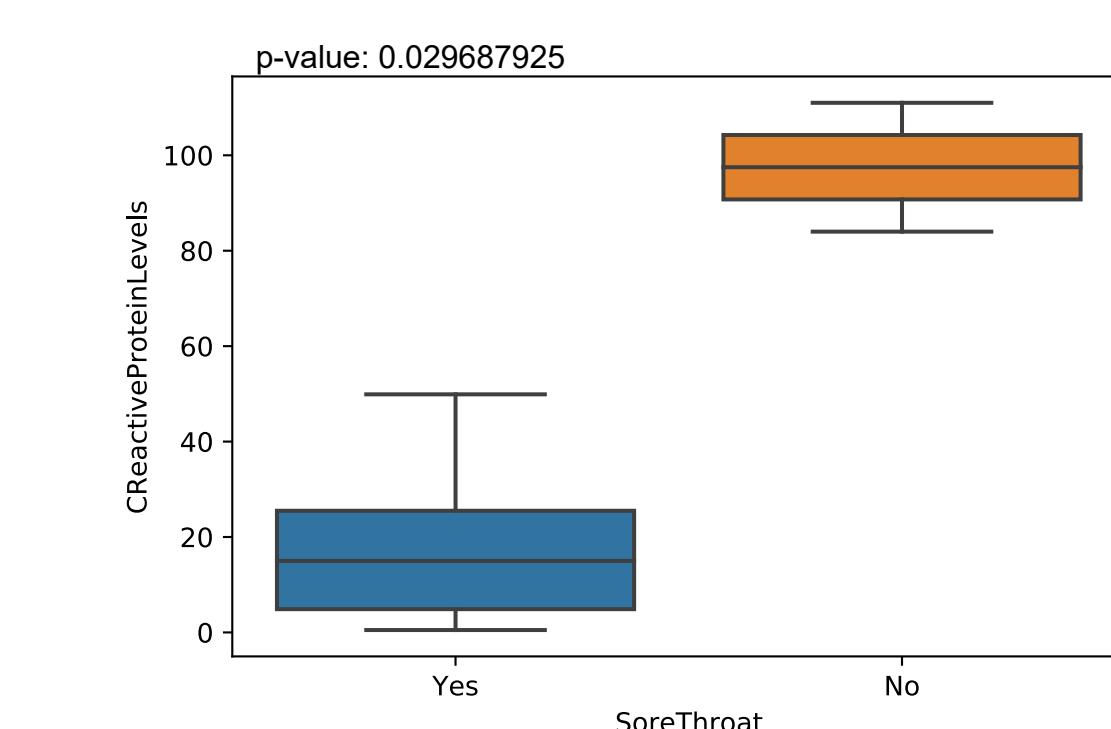
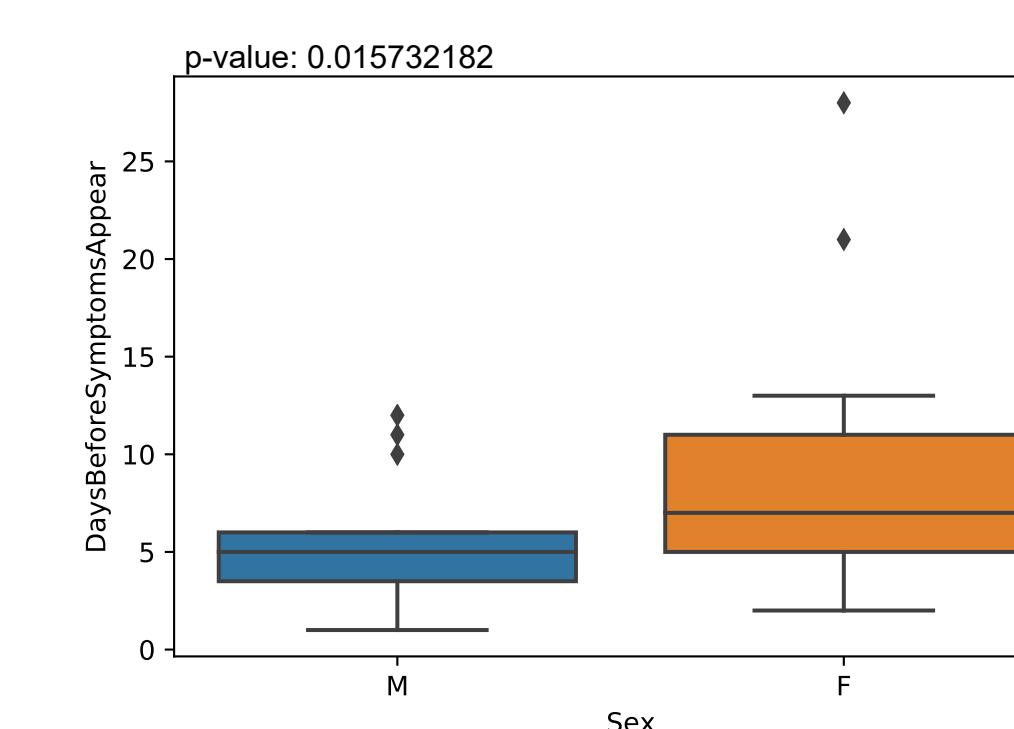
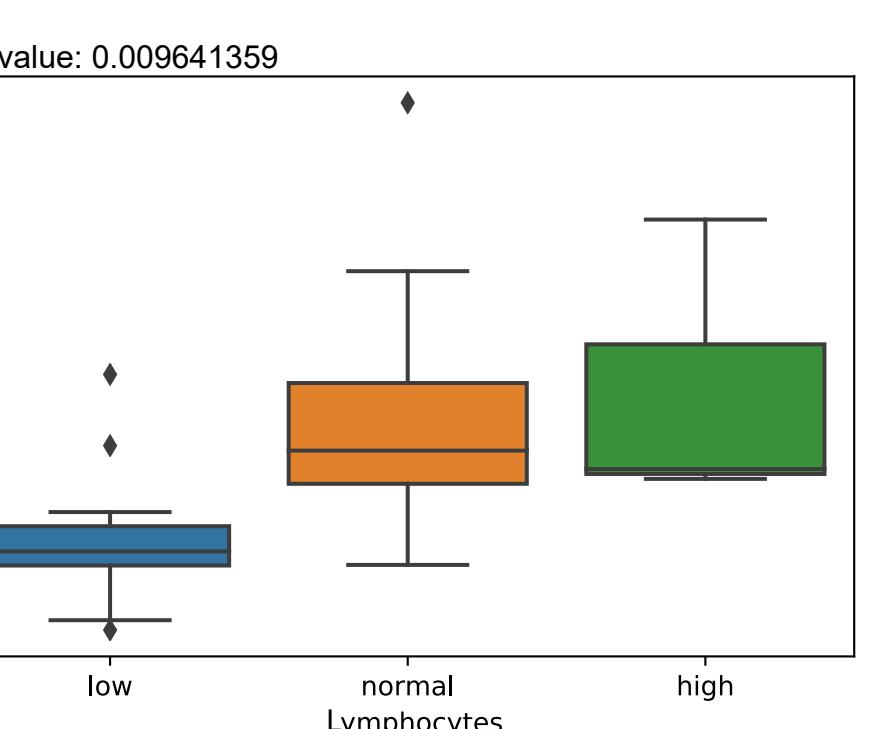
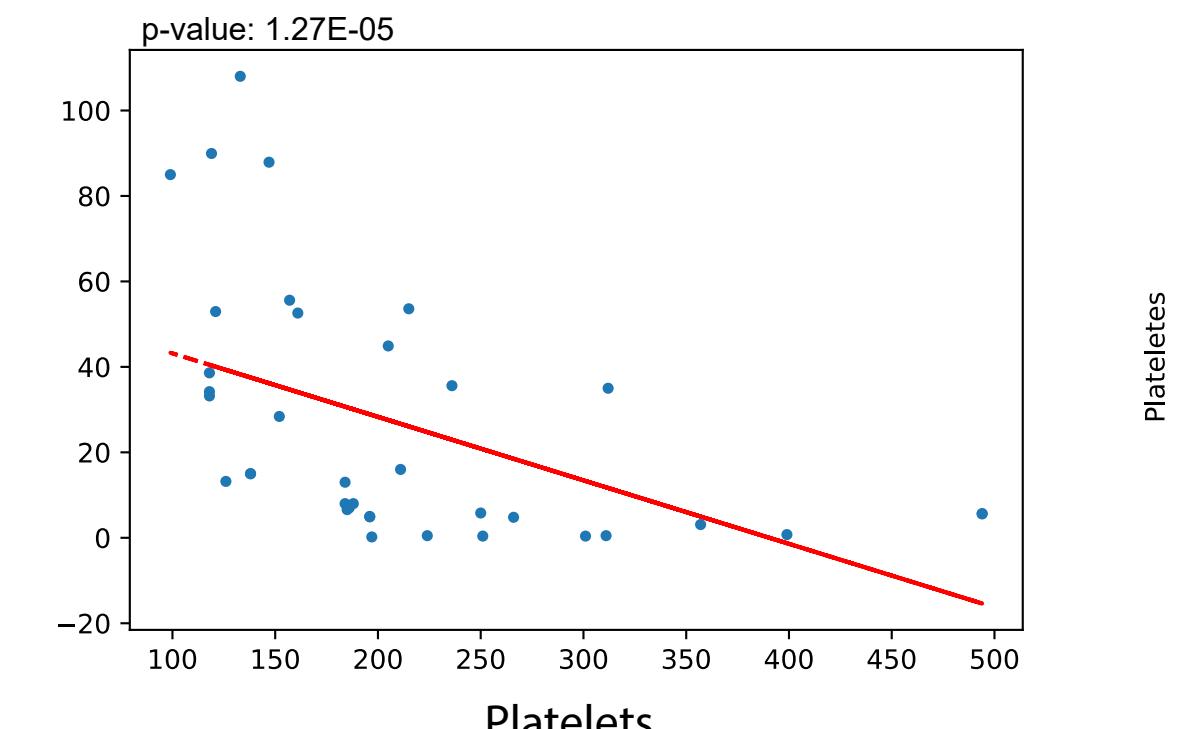
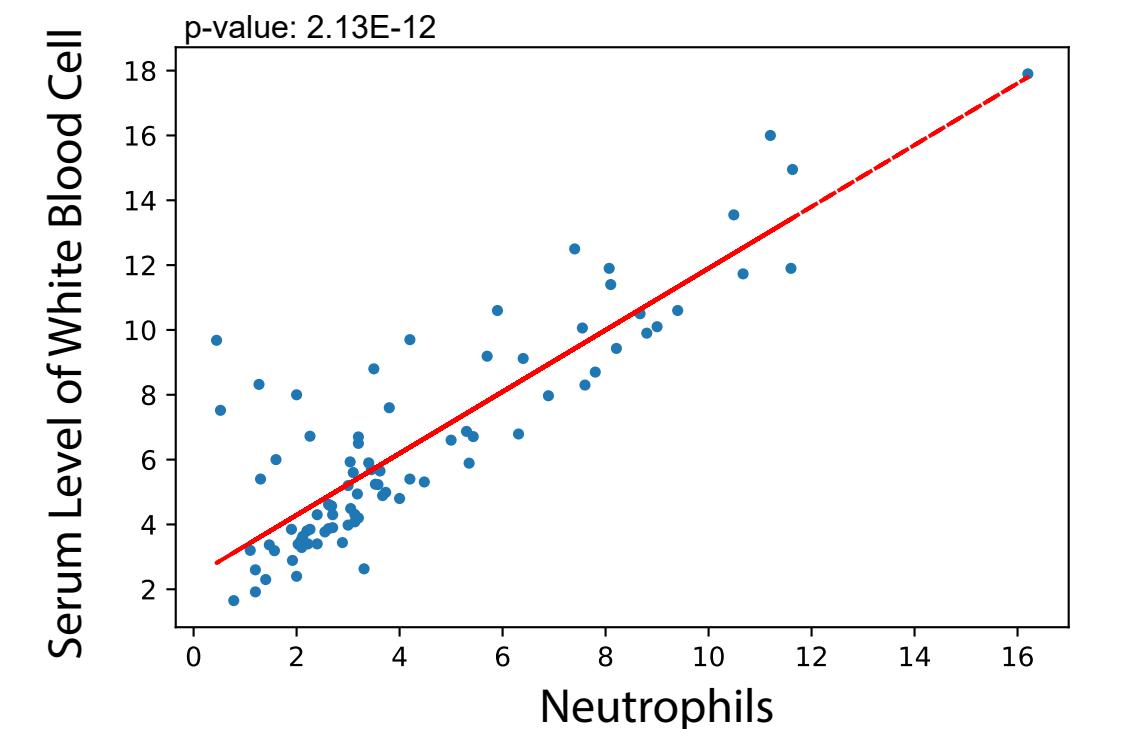
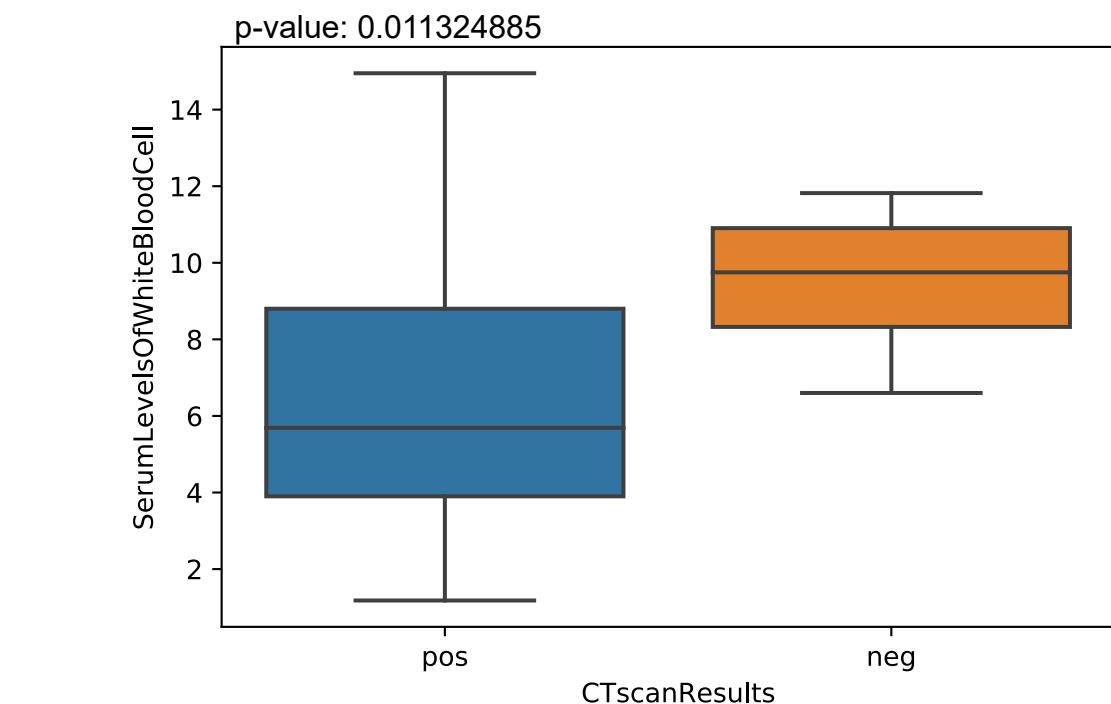
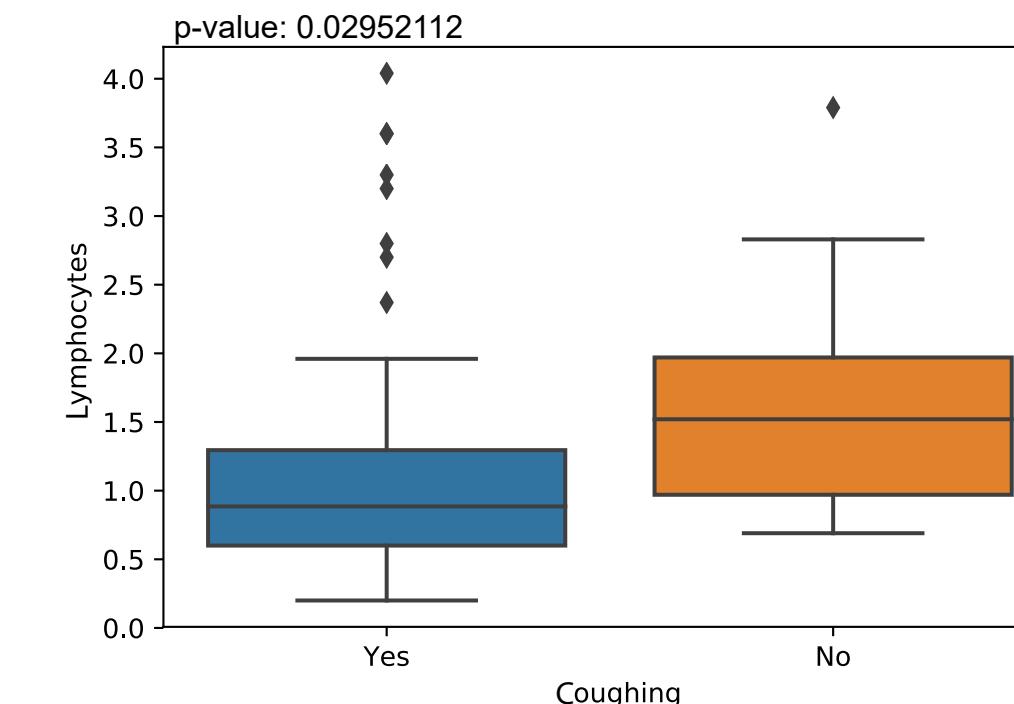
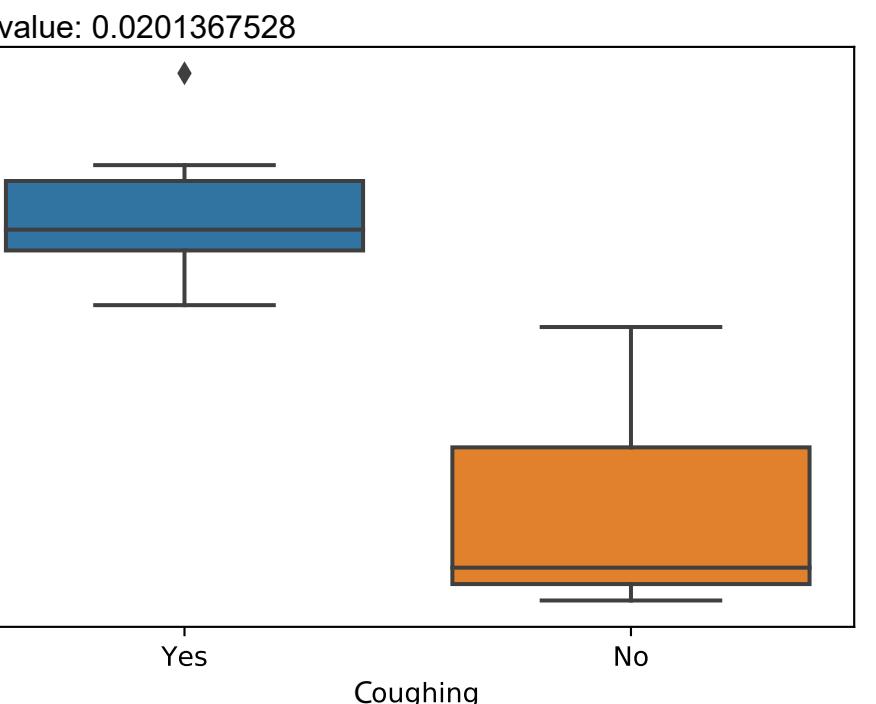
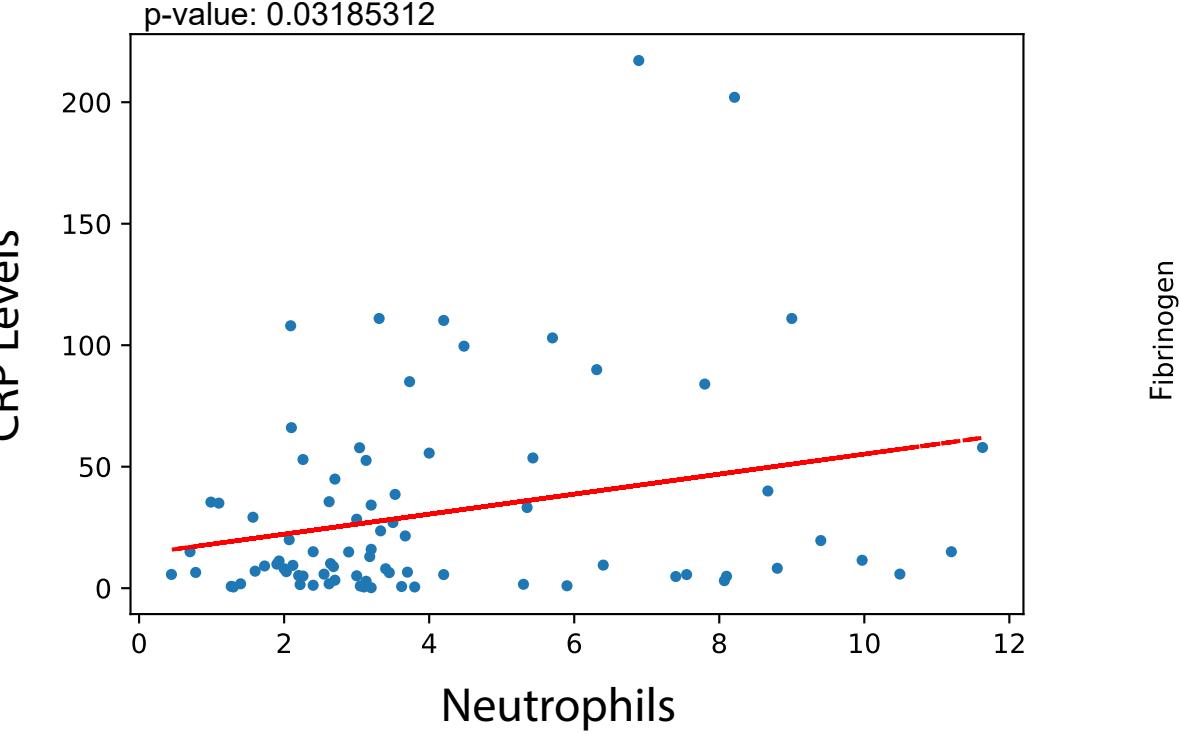
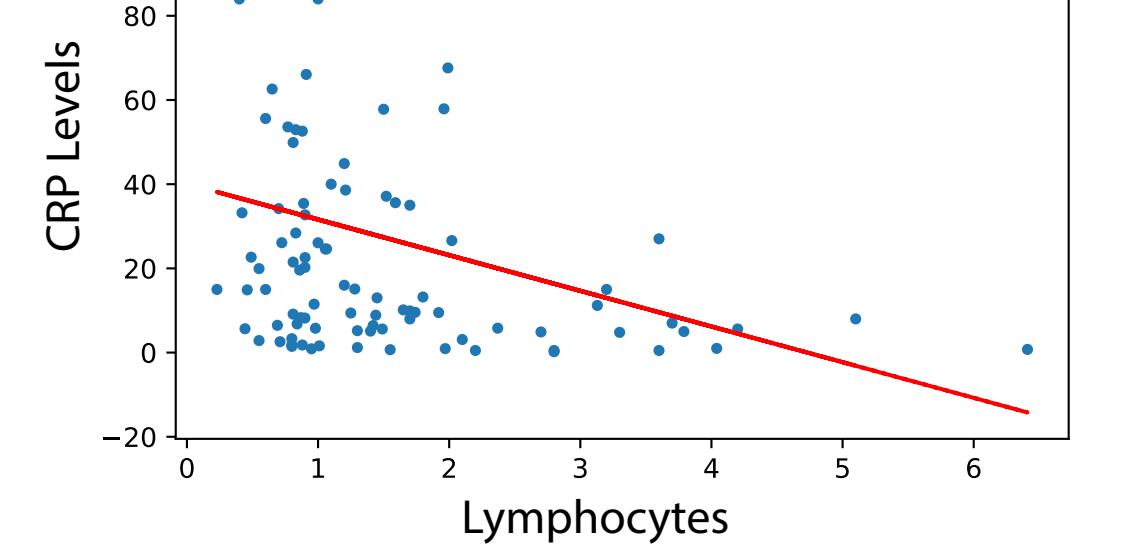
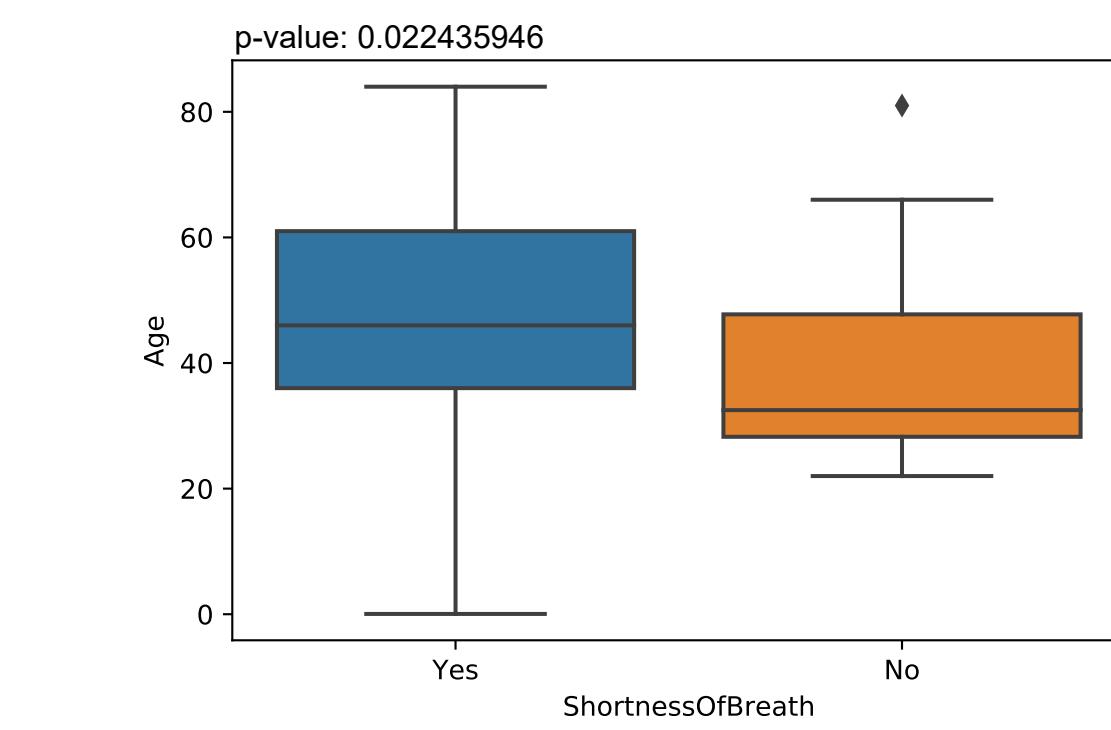
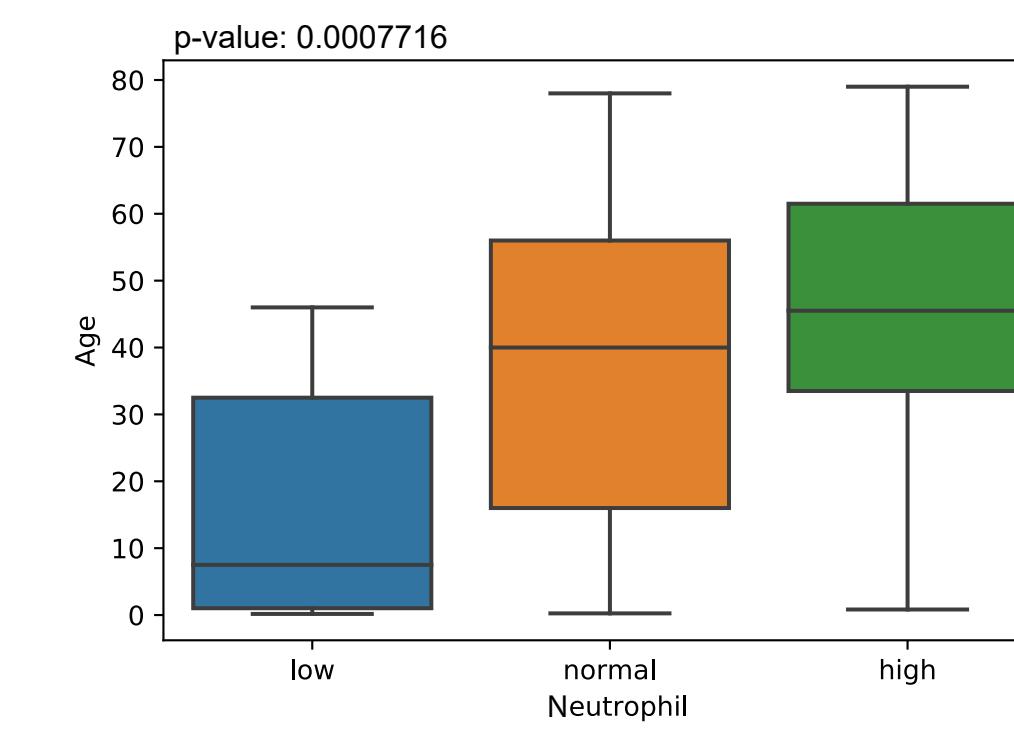
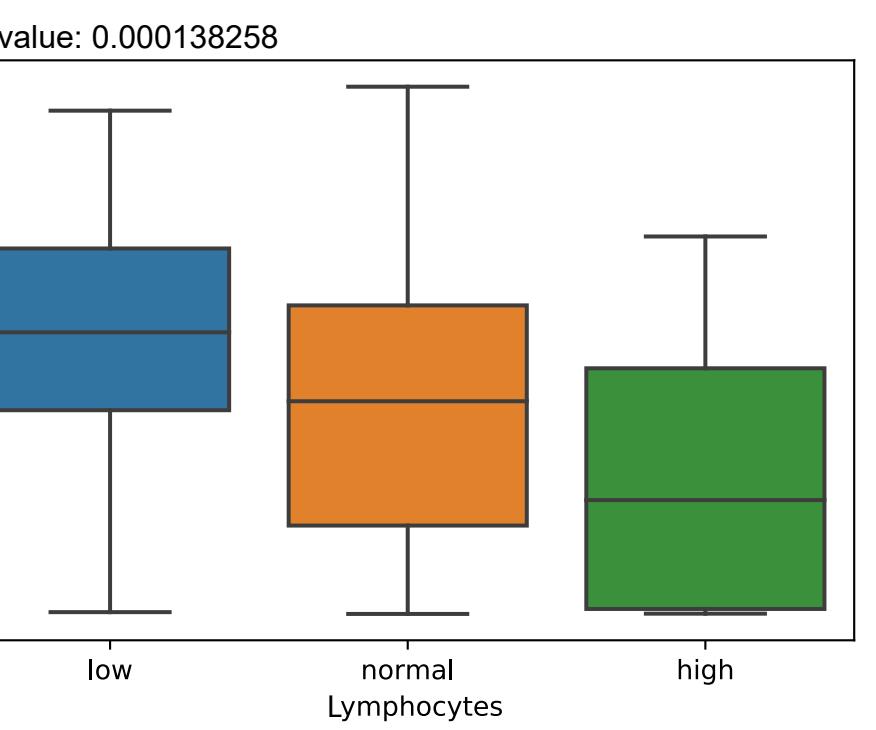
409

410 **Table 1: Clinical Variables Summary of Meta-analysis**

Continuous Variable				
Clinical Variable	# of Data	mean	median	variance
Age	389	38.91306	39	21.85783
NumberOfFamilyMembersInfected	54	3.37037	2	2.6338
neutrophil	103	6.854078	3.31	12.62838
SerumLevelsOfWhiteBloodCell	130	7.031223	5.965	4.250785
lymphocytes	135	2.022841	0.98	4.207139
Plateletes	50	220.32	185.5	146.3334
CRactiveProteinLevels	139	31.18187	15	40.4953
Eosinophils	8	0.06125	0.01	0.070078
RedBloodCells	4	4.225	4.205	0.189011
Hemoglobin	24	45.5	14.5	49.99953
Procalcitonin	33	2.586394	0.07	12.54482
DurationOfIllness	88	14.06818	12	8.970653
DaysToDeath	3	12.66667	12	6.548961
DaysBeforeSymptomsAppear	38	7.368421	6	5.142297
NumberOfAffectedLobes	24	1.75	2	1.163687
TimeBetweenAdmissionAndDiagnosis	47	5.893617	6	4.116568
bodyTemperature	67	37.6209	37.5	0.972999
Hematocrit	7	0.320286	0.355	0.078175
ActivatedPartialThromboplastinTime	9	33.18889	33.4	3.642784
fibrinogen	9	3.685556	3.91	0.752184
urea	19	3.123158	3	0.863884
Discrete Variable				
Variables		Number		Percentage
Sex				
M		194		49.4898
F		198		50.5102
Community Transmission				
Yes		93		37.5

No		46		18.54839
No/Wuhan		109		43.95161
Neutrophil				
low		15		11.81102
normal		83		65.35433
high		29		22.83465
Serum Levels Of White Blood Cell				
low		55		32.35294
normal		94		55.29412
high		21		12.35294
Lymphocytes				
low		86		48.86364
normal		73		41.47727
high		17		9.659091
C Reactive Protein(CRP) Levels				
normal		60		37.97468
high		98		62.02532
CT Scan Results				
pos		124		89.20863
neg		15		10.79137
RT-PCR Results				
pos		100		96.15385
neg		4		3.846154
X-ray Result				
pos		35		74.46809
neg		12		25.53191
GOO				
Yes		92		96.84211
No		3		3.157895
Diarrhea				
Yes		30		45.45455
No		36		54.54545
Fever				
Yes		261		91.25874

No		25		8.741259
Coughing				
Yes		164		82.82828
No		34		17.17172
Shortness Of Breath				
Yes		45		60
No		30		40
Sore Throat				
Yes		37		60.65574
No		24		39.34426
Nausea/Vomiting				
Yes		18		52.94118
No		16		47.05882
Pregnant				
Yes		43		66.15385
No		22		33.84615
Fatigue				
Yes		8		61.53846
No		5		38.46154

A**B**

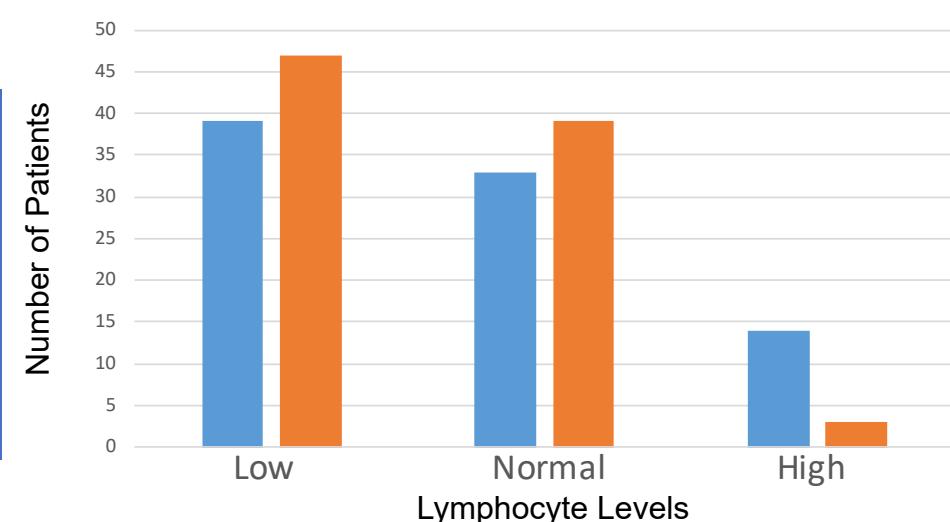
A

Contingency Table: Lymphocyte Levels vs Gender

p-value: 0.02015

Lymphocyte

	Male	Female	All
Low	39	47	86
Normal	33	39	72
High	14	3	17
All	86	89	175



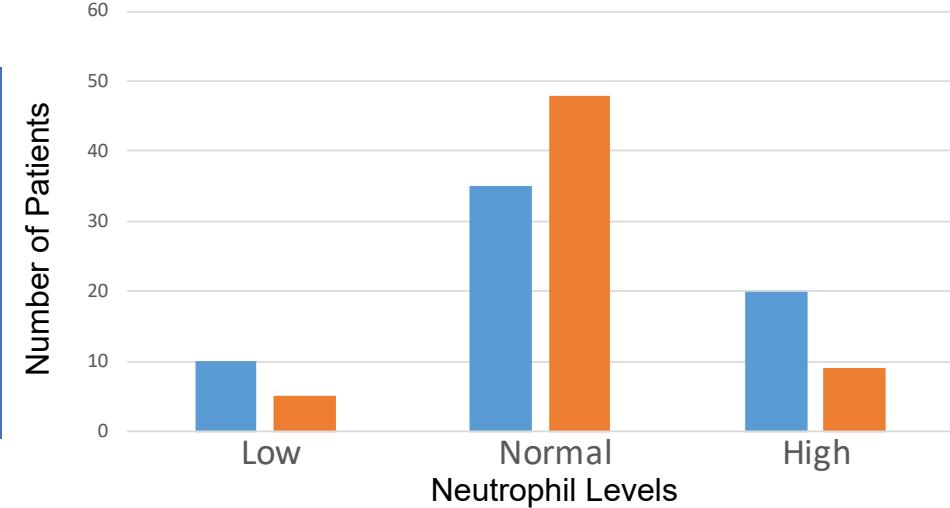
B

Contingency Table: Neutrophil Levels vs Gender

p-value: 0.01566

Neutrophil

	Male	Female	All
Low	10	5	15
Normal	35	48	83
High	20	9	29
All	65	62	127



C

Contingency Table: Serum Leukocyte Levels vs. Gender

p-value: 0.002905

Serum Levels of White Blood Cell

	Male	Female	All
Low	16	39	55
Normal	50	40	90
High	12	7	19
All	78	86	164

