

A Machine Learning Approach to Differentiate Between COVID-19 and Influenza Infection Using Synthetic Infection and Immune Response Data

Suzan Farhang-Sardroodi^{1, 2,*}, Mohammad Sajjad Ghaemi³, Morgan Craig⁴, Hsu Kiang Ooi³, and Jane M Heffernan^{1,2,*}

¹Modelling Infection and Immunity Lab, Mathematics Statistics, York University, Toronto, Canada

²Centre for Disease Modelling (CDM), Mathematics Statistics, York University, Toronto, Canada

³Digital Technologies Research Centre, National Research Council Canada, Toronto, ON, Canada

⁴Sainte-Justine University Hospital Research Centre and Department of Mathematics and Statistics, Université de Montréal, Montreal, Quebec, Canada

*Corresponding Authors: suzanfa@yorku.ca, jmheffer@yorku.ca

January 28, 2022

Abstract

Data analysis is widely used to generate new insights into human disease mechanisms and provide better treatment methods. In this work, we used the mechanistic models of viral infection to generate synthetic data of influenza and COVID-19 patients. We then developed and validated a supervised machine learning model that can distinguish between the two infections. Influenza and COVID-19 are contagious respiratory illnesses that are caused by different pathogenic viruses but appeared with similar initial presentations. While having the same primary signs COVID-19 can produce more severe symptoms, illnesses, and higher mortality. The predictive model performance was externally evaluated by the ROC AUC metric (area under the receiver operating characteristic curve) on 100 virtual patients from each cohort and was able to achieve at least AUC=91% using our multiclass classifier. The current investigation highlighted the ability of machine learning models to accurately identify two different diseases based on major components of viral infection and immune response. The model predicted a dominant role for viral load and productively infected cells through the feature selection process.

Keywords: Biological Systems, Mechanistic Model, Infectious disease, Influenza (flu), COVID-19, Machine Learning, Classification, Logistic Regression, Regularization, Lasso, Ridge

30 1 Introduction

31 Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) and influenza viruses cause
32 COVID-19 and influenza diseases, respectively, and mainly infect the upper and lower respi-
33 ratory tract [1, 2]. Both infections present some similar prime symptoms leading to a clinical
34 dilemma in diagnosing patients with the early infections [3–5]. However, COVID-19 tends to
35 cause worse decompensation due to its intensive transmission, and vascular effects which have
36 led to an unrivaled global crisis [6–9]. Moreover, as the striking COVID-19 outbreak contin-
37 ues, the concurrence of epidemics can be impending. Therefore, it is of interest to design data
38 analysis tools that can accurately differentiate between these two infections and help curb the
39 pandemics.

40 One way to rapidly classify patients as influenza or COVID-19 could be through machine
41 learning approaches. Preliminary investigation illustrated the potentials of machine-learning
42 models for accurately distinguishing between these two viral infections, using demographics,
43 body mass index, and vital signs in infected patients [9]. Herein, we used a simple ML-based
44 classification to identify the patients with influenza and SARS-CoV-2 using mathematically
45 based variables of the in-host infection dynamics and immune response. During the past decade,
46 virus-host mathematical modeling has become an increasingly powerful tool to study intracel-
47 lular viral infection dynamics and the ensuing immune response. Dynamic mathematical mod-
48 eling can deepen our understanding of virus spread within organs that amplify the development
49 of new antiviral drugs, and optimize treatment regimens. Importantly, these models can also
50 help to mitigate difficulties related to clinical data analyses, such as inconsistencies in data col-
51 lection that can lead to biased trial results and significantly complicate comparative analytics.
52 For this purpose, we applied a basic mathematical model on the cellular scale (the so-called
53 target cell-limited model [10, 11]) fit two different sets of *in vivo* data, to create virtual patient
54 cohorts. Using our provided multiclass classifier, the patients were differentiated between the
55 two infections. We certainly hope that our work can be a guide for future applications validated
56 on the external clinical test set to help clinicians and front-line healthcare workers accurately
57 recognize the disease. With just some important in-host measurements clinicians may discern
58 the patients before a laboratory diagnosis.

59 This paper is organized as follows: In section 2, through subsection 2.1, we discuss the
60 In-host mathematical modelling of influenza and COVID-19 and parameter estimation. Via
61 subsection 2.2, we use the mechanistic model to generate synthetic patient data. In subsections
62 2.3 we study developing and evaluating a supervised machine learning method to discriminate
63 the patients with different infections. The Interpretability of the developed model is discussed
64 in subsection 2.4. The results of the prediction are presented in section 3 through subsection
65 3.1. Subsection 3.2 discusses the importance of the data features and determines the dominant
66 features. The paper concludes with a discussion in Section 4.

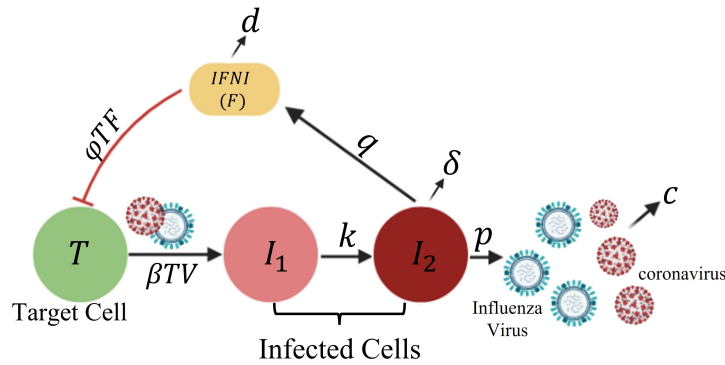


Figure 1: Schematic of viral infection. Each Target cell, T , is infected by a virus, V , with a constant rate β . During the eclipse period the productively infected cell, I_2 , is being produced by the first infected cell, I_1 , with a constant rate k . The Infected cell, I_2 , produces virus at rate p , IFNI at rate q and dies at rate δ per cell. IFNI hinders viral infection by converting target cells to a virus-resistant state with a constant rate ϕ and decays with rate d . Free virus particles that can be influenza or coronaviruses are cleared at per-capita rate c .

67 2 Method

68 2.1 Mechanistic models

We employed a target-cell limited model of viral dynamics using five differential equations that track susceptible target cells (T), infected cells in the eclipse phase (I_1), productively infected cells (I_2), virus (V), and interferon (F) in-host. Figure 1 presents a flow diagram of the model. The system of ordinary differential equations is as follows:

$$\frac{dT}{dt} = -\beta TV - \phi TF \quad (1a)$$

$$\frac{dI_1}{dt} = \beta TV - kI_1 \quad (1b)$$

$$\frac{dI_2}{dt} = kI_1 - \delta I_2 \quad (1c)$$

$$\frac{dV}{dt} = pI_2 - cV \quad (1d)$$

$$\frac{dF}{dt} = qI_2 - dF \quad (1e)$$

69 Briefly, virus particles V can infect susceptible target cells T to produce infected cells. This is
 70 represented by the term βTV . Newly infected cells first enter the eclipse phase I_1 and become
 71 productively infected cells I_2 when within-cell processes that program the cell to make new
 72 virus particles are completed. The eclipse phase takes, on average, $1/k$ time units. Productively

73 infected cells produce new virus particles with a rate of p , and the virus particles are cleared
 74 from the system with a rate of c . We assumed that productively infected target cells have a death
 75 rate δ . Susceptible target cells can be protected from infection by Type I interferon (IFNI),
 76 F . Type I interferons protect neighboring cells from infection and elicit an immune response
 77 [12, 13]. They are central to combating different virus infections and are regularly measured
 78 in clinical trials or infection studies in humans and animals [14]. We assumed that interferon
 79 production is proportional to the number of productively infected cells, that interferon has a
 80 natural decay rate d , and that interferon protects susceptible cells by removing them from the
 81 susceptible target cells population, with a rate ϕF . This term was ignored in [10] for influenza
 82 infection. The model described by Eq. 1 was used in [10] and [11] to examine the kinetics of
 83 influenza A and SARS-CoV-2 viral dynamics, respectively. For the sake of simplicity, we have
 84 ignored a half-day lag in IFNI response that was considered in [10].

85 2.1.1 Parameter Estimation

86 Model parameters for influenza A infection were fit to data from an experimental H1N1 in-
 87 fluenza A/Hong Kong/123/77 infection for six patients [10] and for SARS-CoV-2 from thirteen
 88 untreated patients infected with severe acute respiratory syndrome-coronavirus [11]. The geo-
 89 metric average parameter values along with their 95% confidence intervals and units are sum-
 90 marized in Table 1. We assumed that the initial number of target cells, T_0 , is equal to the total
 91 number of target cells in the upper respiratory tract and set $T_0 = 4 \times 10^8$ cells. In [11] the
 92 authors considered that the target cells distributed in a volume of 30 mL. Assuming that 1%
 93 of these cells expresses the angiotensin-converting enzyme 2 (ACE2) as a receptor for SARS-
 94 CoV-2, the target cell concentration, T_0 , was expressed as 1.33×10^5 cell/mL. Model variables
 with initial values were estimated as in Table 2.

Table 1: Average values and confidence intervals, CI , for influenza A and SARS-CoV-2 within-host viral infection model parameters. Confidence levels of 95% displays the degree of certainty that the parameter values for different samples, fall around the mean.

Influenza Model Parameters [10]									
V_0 [95%CI]	R_0	β [95%CI]	k [95%CI]	p [95%CI]	c [95%CI]	δ [95%CI]	q	d	ϕ
$TCID_{50}/ml^1$		$(TCID_{50}/ml)^{-1}d^{-1}$	d^{-1}	$(TCID_{50}/ml)d^{-1}$	d^{-1}	d^{-1}	d^{-1}	d^{-1}	$d^{-1}cell^{-1}$
0.075[7.6E - 4, 7.5]	21.5[10.1-46.1]	3.2E - 5[6E - 6, 1.7E - 4]	4[3, 5.2]	0.046[0.012, 0.17]	5.2[3.1 - 8.7]	5.2[3.2 - 8.6]	1	1.9 [12, 15, 16]	0
SD:3.5724	SD:17.15	SD:7.8124	SD:1.0486	SD:0.07527	SD:2.6677	SD:2.5724			
COVID-19 Model Parameters [11]									
V_0	R_0 95%[CI]	β	k	p 95%[CI]	c	δ 95%[CI]	q	d	ϕ
$Copies/ml$		$(Copies/ml)^{-1}d^{-1}$	d^{-1}	$(Copies/ml)d^{-1}$	d^{-1}	d^{-1}	d^{-1}	d^{-1}	$d^{-1}cell^{-1}$
0.1	8.6[1.9 - 17.6]	5.68E - 9	3	22.71[0 - 59.64]	10	0.6[0.22 - 0.97]	1	0.4	1.97E-6 [17]
	SD:12.9893			SD:49.3426		SD:0.62051			

¹ $[TCID_{50}/ml]$ corresponds to 4000 $[Copies/ml]$ [18].

² R_0 is the basic reproduction number.

Table 2: Model Variables with Initial values.

Variable	Definition	Initial Value	Unit
T	Target cell	4E+8	Cell
I_1	Infected cell (eclipse phase)	0	Cell
I_2	Productively infected cell	0	Cell
V	Viral load (flu) (COVID-19)	7.5E-2	$TCID_{50}/ml$
		0.1	$Copies/ml$
F	type I interferon (IFNI)	0	Interferon

2.2 Generation of Virtual Patients

To generate a cohort of 100 virtual patients, we followed a technique similar to the one used in [13]. Initial parameter sets representing individual virtual patients were drawn from normal distributions with means fixed to the corresponding parameter value in Table 1 and standard deviations derived from confidence interval measurements. Standard deviations were obtained from standard errors, confidence intervals, and t statistics which measure the size of the difference relative to the variation in the sample data. For each parameter value, the standard deviation was obtained by dividing the length of the confidence interval by standard errors width ($2 \times t - value$) and then multiplying by the square root of the sample size as follows

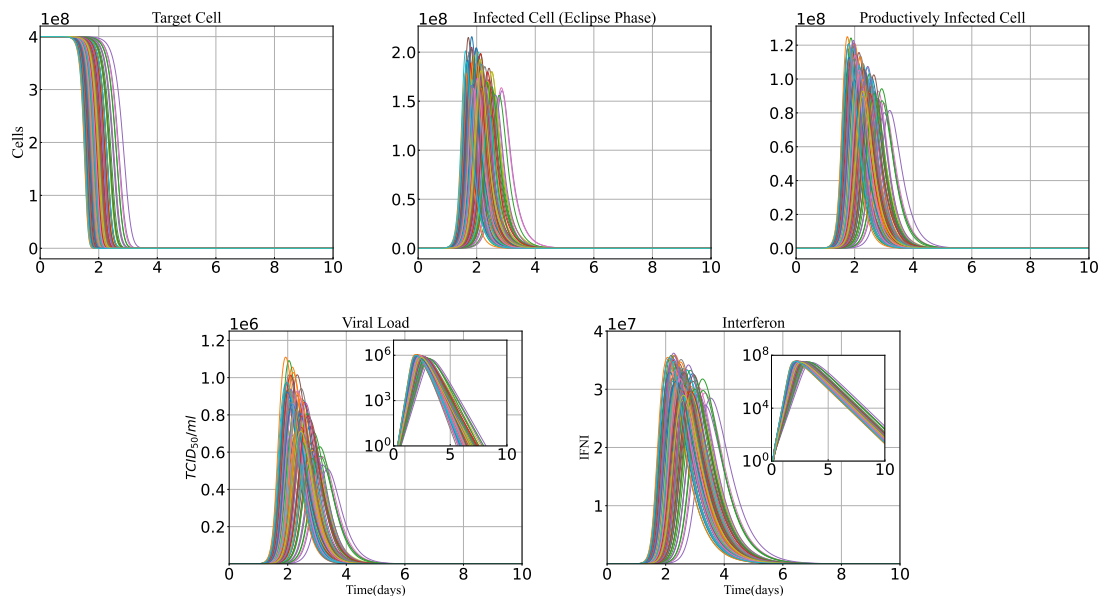
$$SD = \sqrt{N} \times SE = \sqrt{N} \times (upperlimit - lowerlimit)/(2 \times t - value) \quad (2)$$

Standard errors must be of means calculated from within each parameter confidence interval. The $t - value$ for a 95% confidence interval from a sample size of N was then obtained in Microsoft Excel using the $tinv$ function (i.e. $tinv(1 - 0.95, N - 1)$). From [10], the sample size for the influenza cohort is 6 patients infected by H1N1 influenza A/Hong Kong/123/77 infection. The COVID-19 cohort consisted of 13 untreated patients infected with severe acute Respiratory syndrome-coronavirus2 [11]. Therefore, the $t - value$ for influenza patients is 2.571 and for COVID-10 patients is about 2.179. From normal distributions with standard deviations, σ , and means, μ , as the original parameter values, we then generated normal distributions covering values lying around each parameter value such that $|\mu \pm \sigma - \mu| < h$. Herein, the parameter h is the user-defined value as a measure of data diversity. In the other words, the bigger the parameter h , the more diverse the synthetic data. Accordingly, the external noise can affect the data through the parameter h . The dynamics of 100 virtual patients from each cohort are shown in Figure 2. The diversity of patient data is mainly reflected in various viral load levels to agree with prior studies that different viral load is associated with the severity of diseases or different factors such as age or sex of the patients [19].

Consistency of the data

Generating data with time consistency for different cohorts of infections is of great importance. Data inconsistency can lead to loss of information or biased results. Since the influenza mechanistic model predicts faster clearance of influenza-infected cells than SARS-CoV-2 [11], the

Model features for 100 influenza virtual patients



Model features for 100 COVID-19 virtual patients

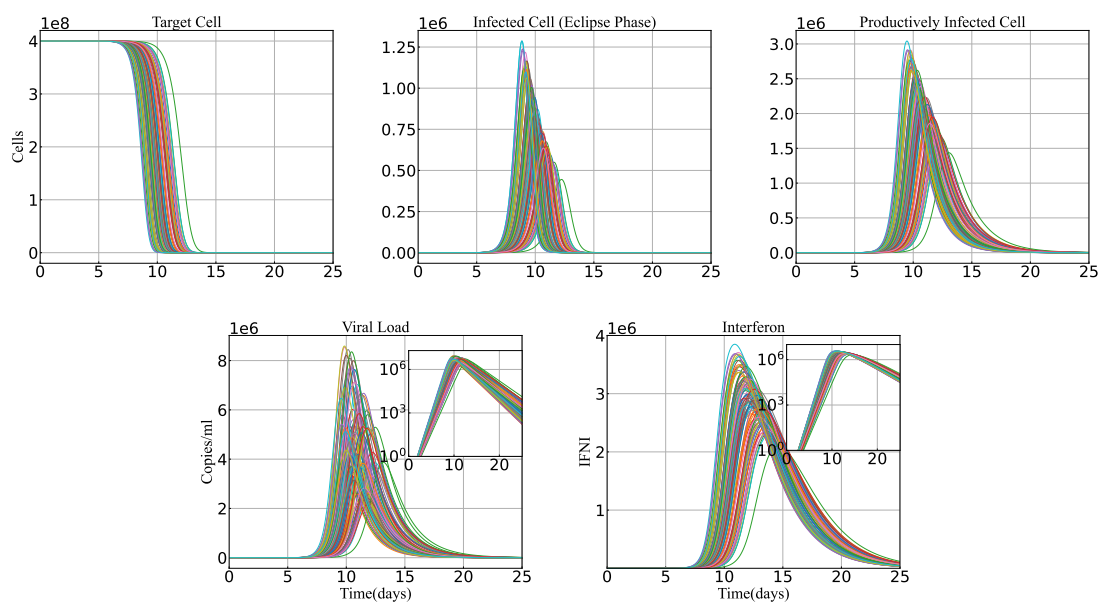


Figure 2: Cohort Dynamics. One hundred virtual patients are generated with different features of Target cells, infected/productively infected cells, viral load, and the only immune factor type I interferon for Influenza (upper two rows) and COVID-19 (lower two rows). Each solid curve with a different color represents a patient. The insets are in log scale.

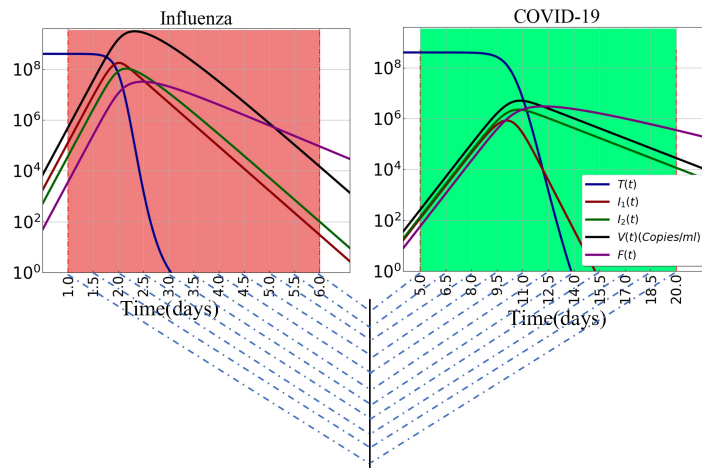


Figure 3: Consistency of the number of virtual data points during the time of infection. Dashed cross blue lines show eleven-time points of an influenza or COVID-19 patient.

125 infection period for influenza and COVID-19 patient dynamics are not the same, see Figure
 126 2. Therefore we limited the consistency of flu/COVID-19 cohorts to have the same number of
 127 data points during the infection time. Hereupon, as an example, we divided the main infection
 128 period (i.e., $[1 - 6]$ days for influenza patients and $[10 - 20]$ days for COVID-19 patients)
 129 into ten different sub-intervals with half-day length time steps for influenza patients and
 130 half-day length time steps for COVID-19 patients (see Figure 3). Hence, despite having differ-
 131 ent infection periods and time steps with different lengths to report the new virtual data point,
 132 the total number of data for the two different cohorts was the same.

133 In addition to the total infection period, we were also interested in studying the viral load
 134 dynamics in the early period of infection. The median incubation period for influenza A(B)
 135 virus is estimated to be 1.4(0.6) days, and for SARS-CoV-2 is around 5–6 days [20]. Therefore,
 136 we assumed the time interval $[0.9, 1.3]$ days for influenza, and $[5 - 6.5]$ days for COVID-19
 137 cohorts, corresponding to $[10^2 - 10^4]$ *Copies/ml* viral load. Dividing each interval into three
 138 different sub-intervals to get the time steps with length one-sixth of a day for Influenza and
 139 half a day for COVID-19 patients, we had four consistent data points for each patient.

140 2.3 Predictive Model Development

141 To distinguish between patients who encounter COVID-19 from those who are exposed to in-
 142 fluenza, we developed a predictive model based on some biological feature selections. Accord-
 143 ingly, we adopted Logistic regression with ℓ_1 -regularization, referred to Lasso (stands for least
 144 absolute shrinkage and selection operator) Regression, as an appropriate technical classifica-
 145 tion. Lasso regression is widely used for many supervised classification problems based on
 146 the concept of probability [21]. It can simplify the model complexity by removing irrelevant

147 features of the data set. Recently, this algorithm was used by Han et al. to find some addi-
148 tional novel immune features that accurately identified patients before the clinical diagnosis of
149 preeclampsia [22].

150 Logistic regression, which is a special case of linear regression and used for binary classifi-
151 cation, is defined by the following sigmoid function

$$152 \quad h(X) = \frac{1}{1 + e^{-(\beta_0 + \beta X)}} \quad (3)$$

153 in which X is the $(n \times p)$ model feature matrix of $n = 100$ patients and $p = 5$ biological
154 hallmarks. Defining the cost/objective (C) function of logistic regression in mean squared error
155 format leads to a non-convexity that makes it difficult to optimally converge. Therefore, it is
156 represented by the following equations

$$157 \quad C(h(X), Y) = \begin{cases} -\log(h(X)), & \text{if } y = 1 \\ -\log(1 - h(X)), & \text{if } y = 0 \end{cases} \quad (4)$$

158 where Y is a binary response vector of outcome (CVOID-19 vs flu). Compressing the above
159 two equations inside a single function, we have

$$160 \quad J(X) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(h(x_i)) + (1 - y_i) \log(1 - h(x_i))] \quad (5)$$

161 Replacing the sigmoid function from equation (3) and applying a penalty term equal to the
162 absolute value of the magnitude of coefficients, we can reach the following objective function
163 (after doing some mathematical simplifications) [22]

$$164 \quad J(X) = -\left[\frac{1}{n} \sum_{i=1}^n y_i (\beta_0 + \beta^T x_i) - \log(1 + e^{\beta_0 + \beta^T x_i}) \right] + \alpha \|\beta\|, \quad \alpha > 0 \quad (6)$$

165 The penalty term which is called the ℓ_1 -regularization term is added to prevent data over-fitting.
166 The model objective is to find a specific solution with a best-optimized cost function.

167 For model training and testing, we developed a \mathcal{K} -fold cross-validation strategy, which is
168 a re-sampling method to evaluate machine learning models on a limited data sample. The
169 procedure has a single parameter called \mathcal{K} which displays the number of groups that a given
170 data sample is to be split into. As such, the procedure is often called \mathcal{K} -fold cross-validation.
171 Therefore, our regression model is not tailored to a particular data set and is exposed to all
172 available samples of a given subject in the training set. This approach implies that the training
173 procedure was entirely blinded to the synthetic patient data sets, and ensures the presumed
174 independence from any intra-subject correlations that are required for Lasso classification. We
175 fixed the number of folds of the data as $\mathcal{K} = 5$. Running the analysis on each fold, the predicted
176 outcome will be the one with the least estimated prediction error. The regularization parameter
177 α is estimated by a cross-validation procedure.

178 **2.3.1 Evaluating Model Performance**

179 The discriminating ability of the developed model in predicting patients with influenza from
180 COVID-19 was evaluated using AUC (Area Under The Curve) ROC (Receiver Operating Char-
181 acteristics) curve analysis. AUC - ROC curve is one of the most important evaluation metrics to
182 visualize the performance of multi-class classification problems. ROC represents a probability
183 curve of sensitivity (true positive rate= $\frac{TP}{TP+FN}$) against 1-specificity (false positive rate= $\frac{FP}{FP+TN}$)
184 and AUC is a performance measure of discrimination. In the other words, the AUC score is a
185 criterion that explains how well the model is capable of discerning different cohorts. Generally,
186 an AUC closer to 1 indicates a better overall diagnostic performance of influenza classes as
187 influenza or COVID-19 to COVID-19.

188 **2.4 Model Interpretability**

189 From [23, 24], "Interpretability" is the degree to which a human can understand the cause of
190 a decision and consistently predict the model's result. The higher the interpretability of a ma-
191 chine learning model, the better understanding of why certain predictions have been made.
192 Interpretable machine learning models are beneficial to extract the relevant knowledge from
193 relationships either contained in data or learned by the model [25, 26].

194 Here, we looked at the regularization path which is a plot of all coefficients values against
195 the values of α in ℓ_1 penalization term, to see the behavior of the Lasso regression and interpret
196 the prediction outcomes. The main purpose of Lasso regression is to classify groups of data
197 by providing feature coefficients that can select the important features and maintain model reg-
198 ularization to avoid over-fitting the data. Therefore, the Lasso path can give us an idea of the
199 feature's importance.

200 **3 Results**

201 **3.1 Prediction of Influenza versus COVID-19 infection**

202 In this study we developed a classifier in the Lasso framework to identify patients with either in-
203 fluenza or COVID-19, based on four major entities of viral dynamics, $\{T(t), I_1(t), I_2(t), V(t)\}$,
204 and one main factor of host immune response, type I interferon ($F(t)$), as the entry data fea-
205 tures. The model was trained on data from one hundred virtual patient-level data in each infec-
206 tion cohort without noise, and it was externally validated on testing set with demographic noise
207 (reflected in diverse viral load levels). Results in Figures 4, 5 and 6 reflect the Lasso predictions
208 using the entire infection period (see Section 2.2). In Figures 4 and 5, two-dimensional scatter
209 plots are used to compare ground truth to regression predicted values based on all model fea-
210 tures. The hue spectrum from light to dark illustrates the probability of being in the influenza
211 (blue) or COVID-19 (red) group. In the other words, the darker the colors, the better the pre-
212 diction. Considering three attributes in the data, the predicted outcomes are improved. This

213 is shown in three-dimensional scatter plots in Figure 6 of the ground truth and regression pre-
214 dicted values. ROC AUC=95% indicates a satisfactory performance of the model to distinguish
215 between COVID-19 and influenza patients, see figure 7 for more details.

216 **Early days of infection**

217 We examined the model prediction for the data generated at the early days of infection after
218 the incubation period. The results are shown in Figure 8 based on the model features. From
219 the figure, we can see that there are some mispredictions, for small values of $I_1(t)$, $I_2(t)$, $V(t)$,
220 and $F(t)$, especially when $I_2(t)$ is plotted as a function of $I_1(t)$ or $V(t)$ is plotted in terms of
221 $I_2(t)$. In the other words, for this range of values, the influenza patients were misdiagnosed with
222 COVID-19. In an attempt to find the reason, we compared correlations between the different
223 variables in our model – see Figure 9). Here, we see small regions of overlap between influenza
224 and COVID-19 models. Accordingly, the compatibility of the results between the two infections
225 may lead to some overlaps in the model predictions. However, the ability of the model in the
226 prediction of infections when the patients were monitored by $V(t)/F(t)$ as a function of $I_1(t)$,
227 panels (b) and (c), or $F(t)$ in terms of $I_2(t)/V(t)$, panels (e) and (f), can be satisfactory, and
228 thus can serve as benchmarks for clinical diagnosis. The model had a ROC AUC of 91% on the
229 external validation data set for early infection – see Figure 7.

230 **3.2 Significance of the features**

231 To investigate the importance of various data features we created our ℓ_1 -regularization path,
232 which was the best way to see the behavior of the Lasso regression. The regularization path
233 is a plot of all coefficient values in terms of the regularization parameter. Figure 10 illustrates
234 the selection path of each feature with its corresponding coefficient in terms of the logarithm
235 of the regularization parameter α . For each value of α , the path method on the Lasso object
236 returns the coefficients that solve the logistic regression problem with that parameter value. The
237 optimal value of $-\log(\alpha)$ was estimated at around 3.25 for the test set distributed over the entire
238 infection course, and 3.04 when the early days of infection were studied. The results suggested
239 a higher coefficient value for viral load $V(t)$ and productively infected cells $I_2(t)$ compared to
240 the other features.

241 **4 Discussion**

242 This study presents a machine learning model to effectively classify influenza and COVID-19
243 virtual patients using in-host patient data. Our model employed a Lasso regression classifier
244 trained to identify between two hundred patients, highlighted by a ROC AUC of 95%. Using
245 the existing within-host models, We generated synthetic data with five in-host measurements
246 including target cells, eclipse phase, and productively infected cells, viral load, and type I IFN.
247 Analyzing the feature importance revealed that the viral load and the productively infected cells

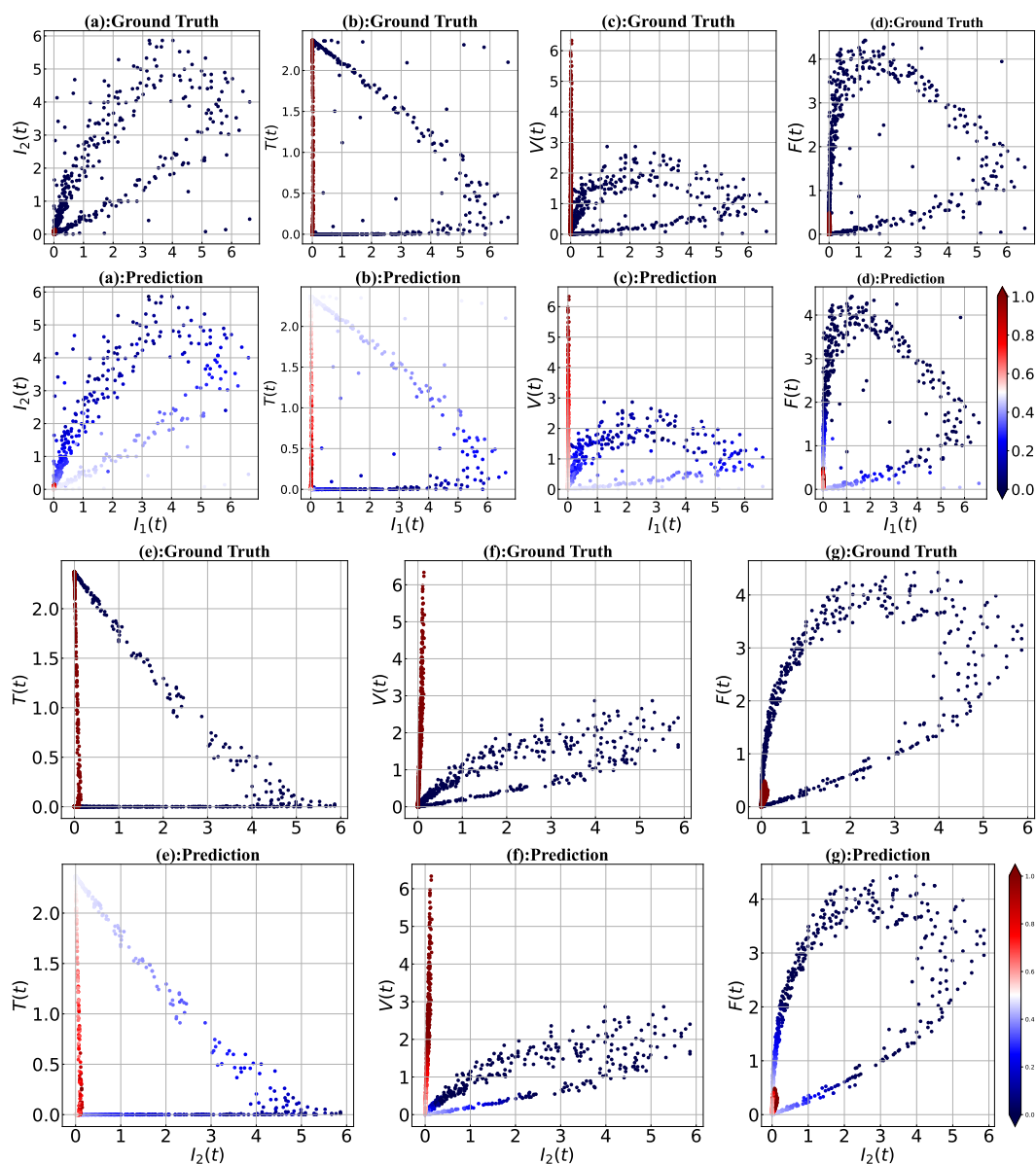


Figure 4: Two-dimensional scatter plots of ground truth and regression predicted values based on model features. Classification of the data was done for: I_2 versus I_1 in panels (a), T vs. I_1 in panels (b), V vs. I_1 in panels (c), F vs. I_1 in panels (d), T vs. I_2 in panels (e), V versus I_2 in panels (f), and F versus I_2 in panel (g). Color denotes the patient probability of being in the influenza (blue color scheme) or COVID-19 (red color scheme) cohorts. Data points, corresponding to each model feature, are rescaled by dividing by their standard deviations.

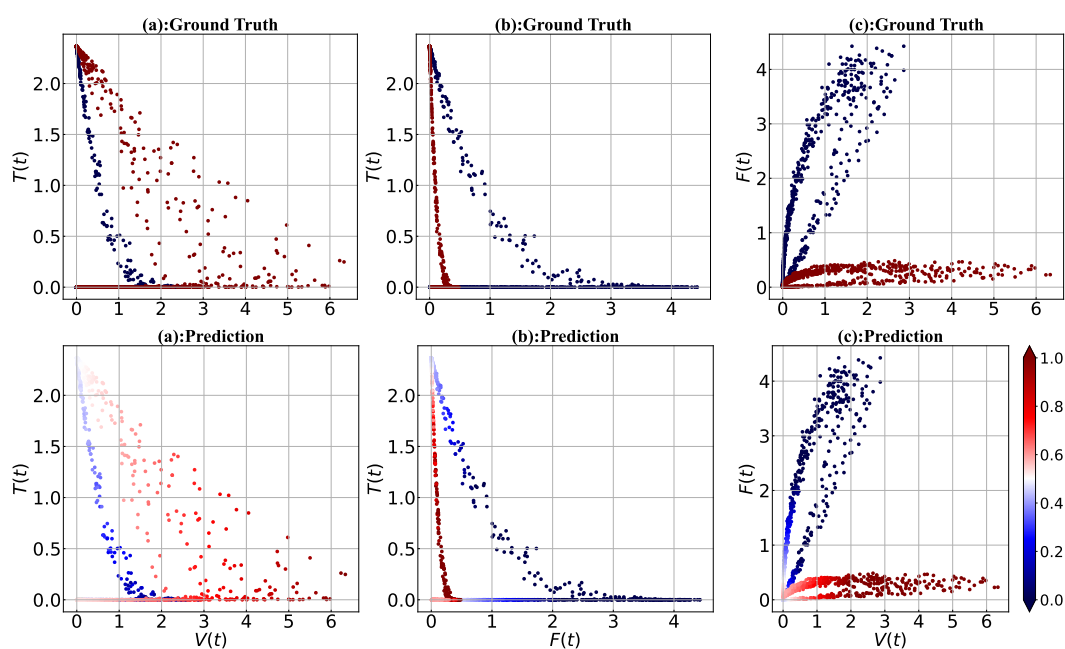


Figure 5: Two-dimensional scatter plots of the ground truth and regression predicted values for three model features T , V , F . Classification of the data was done based on: T versus V in panels (a), T versus F in panels (b), and F versus V in panels (c). Color denotes the patient probability of being in the influenza (blue color scheme) or COVID-19 (Red color scheme) cohorts.

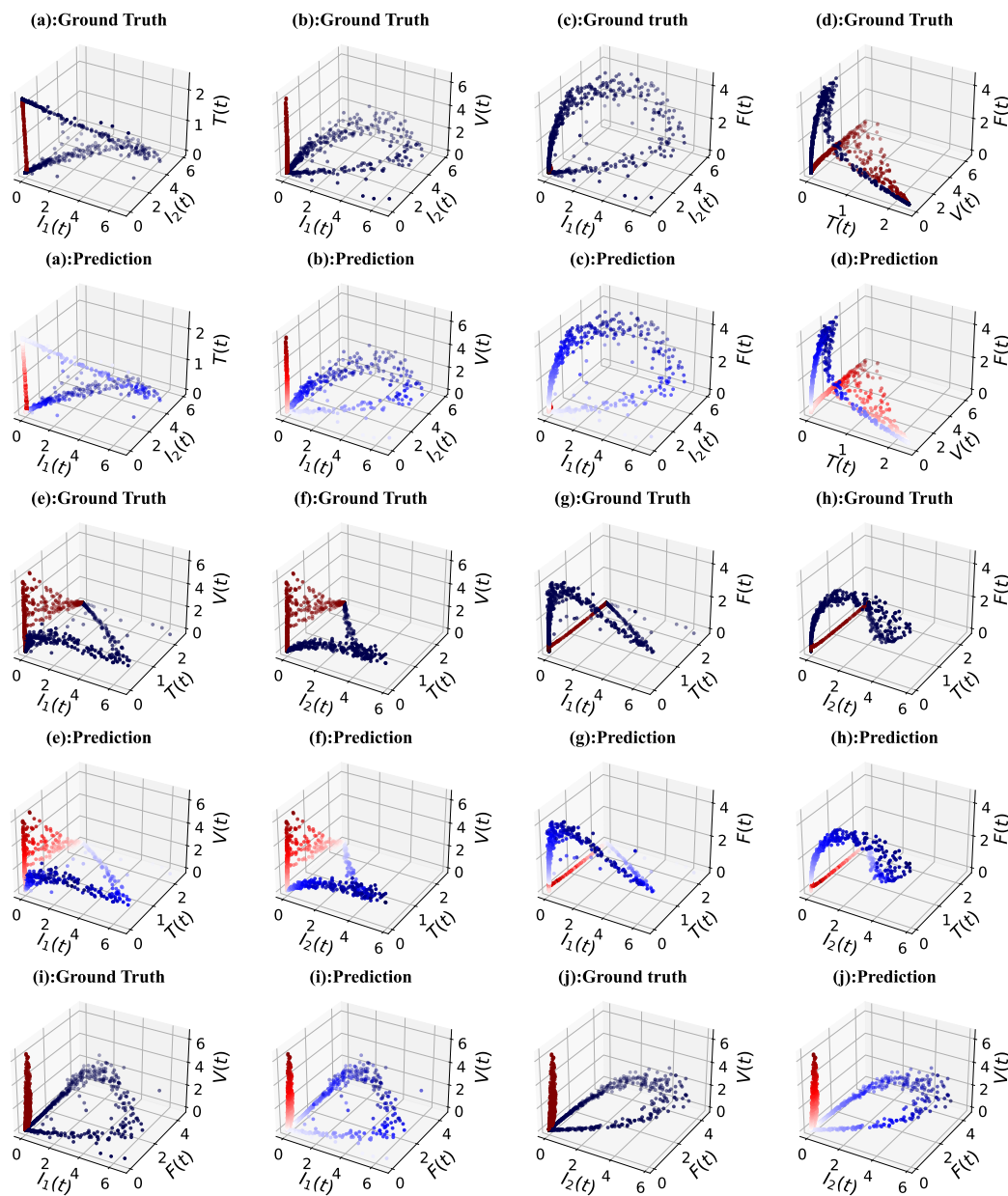


Figure 6: Three-dimensional scatter plots of the ground truth and regression predicted values based on all model features. Classification is based on I_1, I_2, T in panels (a), I_1, I_2, V in panels (b), I_1, I_2, F in panels (c), T, V, F in panels (d), I_1, T, V in panels (e), I_2, T, V in panels (f), I_1, T, F in panels (g), I_2, T, F in panels (h), I_1, F, V in panels (i), and I_2, F, V in panels (j). Shades of blue (red) indicate influenza (COVID-19) group patients. Data points are dimensionless by dividing by the corresponding standard deviations.

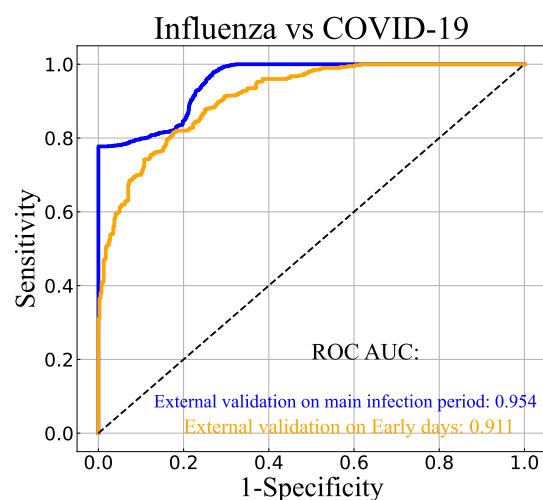


Figure 7: Receiver Operating Characteristic curve (ROC) of influenza vs COVID-19 patients. Area under the ROC curve indicates the predictive performance of the model between COVID-19 and influenza encounters on the external validation test during the main infection (blue curve) and early days (orange curve). The black dashed line in the diagonal has a ROC AUC of 0.5.

248 are the most important components to determine if a patient is infected by influenza or SARS-
249 CoV-2.

250 While our machine learning model was built on the synthetic data distributed during the
251 main infection period, it ascertained a good performance (ROC AUC = 91%) even for the early
252 days, once after the incubation period. However, there are some exceptions for the small values
253 of in-host features where the influenza patients are misdiagnosed by COVID-19 for the early
254 days of infection. The reason was explained by the fact that during the early days of the in-
255 fection, influenza and COVID-19 patients have comparable in-host measurements that lead to
256 some errors in discriminating the patients. This is interpreted as a limitation of our model and
257 can be a future extension of developing dynamic models which take more immune entities into
258 account and end in a better classifier.

259 Our model was trained and successfully evaluated on synthetic data. The model, however,
260 could be applied to animal or human clinical data. This could be useful, for example, if a clinical
261 trial is complicated by the existence of an infectious disease with similar infection character-
262 istics. The model could be applied as a low-cost classification system that would not require
263 expensive virus typing procedures and could rely solely on viral load and interferon measure-
264 ments. We note that studies like [9] that focus analysis on demographic and observational data
265 can be cheaper to conduct, but these data can also be subject to inconsistencies and bias, af-
266 fecting classification outcomes. In a future study, we will expand our analysis to a model of
267 in-host measurements and observational data to determine if specific combinations of in-host
268 and observational data that best classify influenza and COVID-19 infections differ.

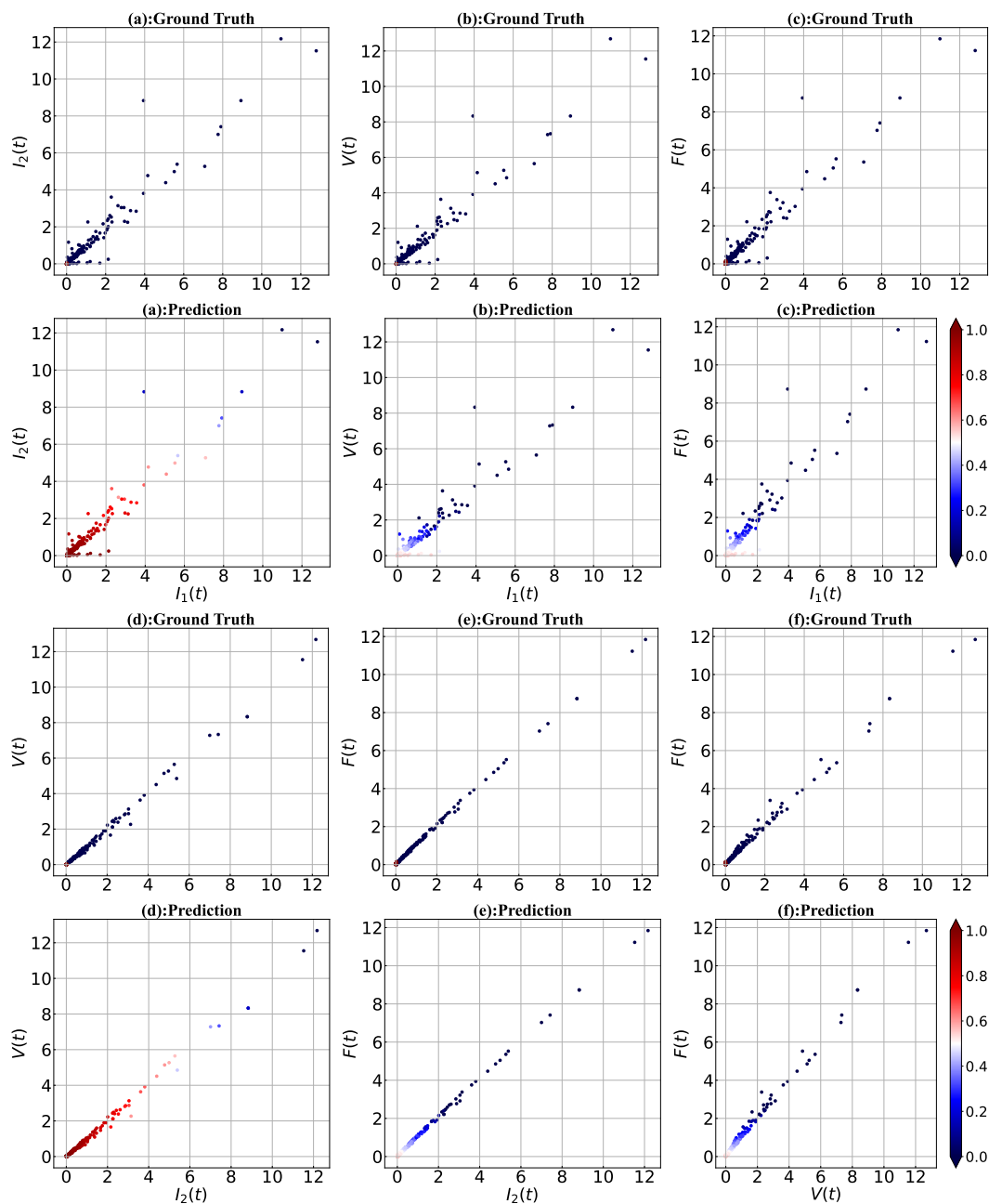


Figure 8: Early days of infection. Two-dimensional scatter plots of the ground truth and regression predicted values based on model features are shown. Classification is based on I_1, I_2 in panels (a), I_1, V in panels (b), I_1, F in panels (c), I_2, V in panels (d), I_2, F in panels (e), and V, F in panels (f). Shades of blue (red) indicate influenza (COVID-19) group patients. Data points are dimensionless by dividing by the corresponding standard deviations.

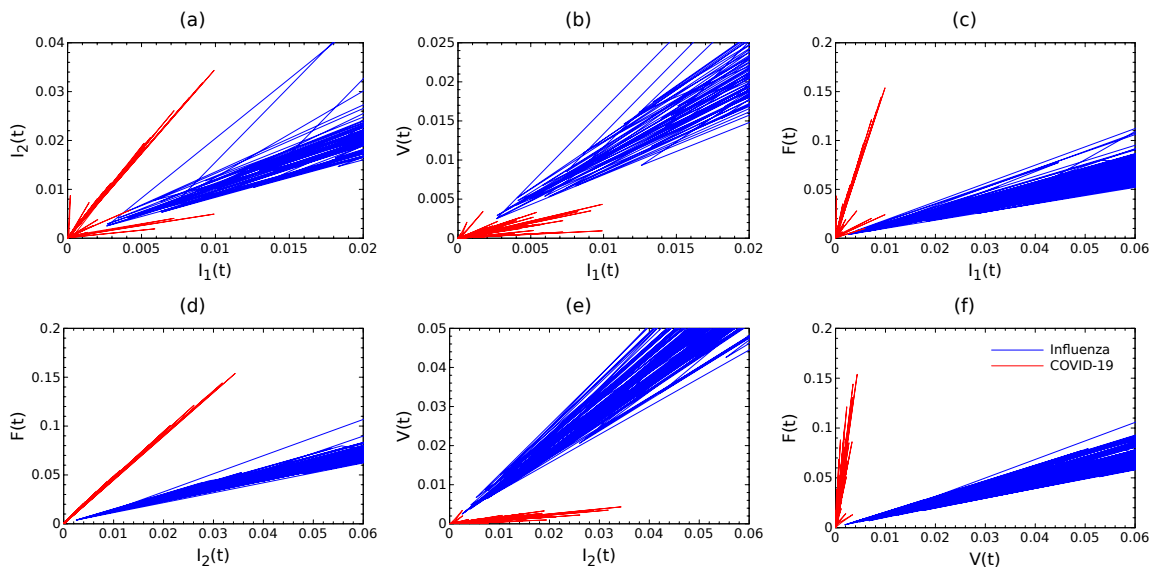


Figure 9: Comparison of in-host measurements, $\{T, I_1, I_2, V, F\}$, between influenza and COVID-19 virtual patients where plotted as a function of each other. Blue(red) solid lines represent the ratio of the features for one hundred influenza (COVID-19) patients. Data points are divided by the corresponding standard deviations for each feature.

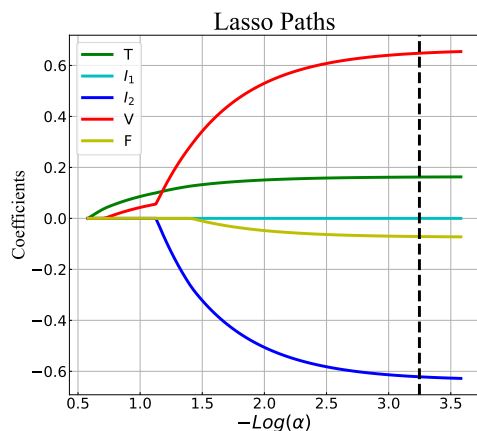


Figure 10: Lasso coefficients of five sample features, $\{T, I_1, I_2, V, F\}$, as a function of the logarithm of regularization parameter, $-\log \alpha$. Each colored line represents the value taken by a different coefficient in the optimization objective for Lasso. The black dashed line indicates the selected regularization parameter with the value of $-\log(\alpha) \approx 3.25$. This number was ≈ 3.04 with the same Lasso Paths when the early days of the infection period were considered.

269 Our machine learning model was developed in the Lasso framework. Ridge regression could
270 also be employed, and require only small changes to our method to include this. We find that
271 the model demonstrated a satisfactory performance by using a Ridge regression classifier –
272 (ROC AUC= 95%) for the main infection period, and (ROC AUC= 89%) for the early days of
273 infection.

274 References

- 275 [1] Kanta Subbarao and Siddhartha Mahanty. Respiratory virus infections: understanding
276 covid-19. Immunity, 52(6):905–909, 2020.
- 277 [2] Laura D Manzanares-Meza and Oscar Medina-Contreras. Sars-cov-2 and influenza: a
278 comparative overview and treatment implications. Boletín médico del Hospital Infantil de
279 México, 77(5):262–273, 2020.
- 280 [3] Shuhei Azekawa, Ho Namkoong, Keiko Mitamura, Yoshihiro Kawaoka, and Fumitake
281 Saito. Co-infection with sars-cov-2 and influenza a virus. IDCases, 20:e00775, 2020.
- 282 [4] Pavan K Bhatraju, Bijan J Ghassemieh, Michelle Nichols, Richard Kim, Keith R Jerome,
283 Arun K Nalla, Alexander L Greninger, Sudhakar Pipavath, Mark M Wurfel, Laura Evans,
284 et al. Covid-19 in critically ill patients in the seattle region—case series. New England
285 Journal of Medicine, 382(21):2012–2022, 2020.
- 286 [5] Hossein Khorramdelazad, Mohammad Hossein Kazemi, Alireza Najafi, Maryam
287 Keykhaee, Reza Zolfaghari Emameh, and Reza Falak. Immunopathological similar-
288 ities between covid-19 and influenza: Investigating the consequences of co-infection.
289 Microbial pathogenesis, 152:104554, 2021.
- 290 [6] Maximilian Ackermann, Stijn E Verleden, Mark Kuehnel, Axel Haverich, Tobias Welte,
291 Florian Laenger, Arno Vanstapel, Christopher Werlein, Helge Stark, Alexandar Tzankov,
292 et al. Pulmonary vascular endothelialitis, thrombosis, and angiogenesis in covid-19. New
293 England Journal of Medicine, 383(2):120–128, 2020.
- 294 [7] Qun Li, Xuhua Guan, Peng Wu, Xiaoye Wang, Lei Zhou, Yeqing Tong, Ruiqi Ren,
295 Kathy SM Leung, Eric HY Lau, Jessica Y Wong, et al. Early transmission dynamics
296 in wuhan, china, of novel coronavirus–infected pneumonia. New England journal of
297 medicine, 2020.
- 298 [8] Na Zhu, Dingyu Zhang, Wenling Wang, Xingwang Li, Bo Yang, Jingdong Song, Xiang
299 Zhao, Baoying Huang, Weifeng Shi, Roujian Lu, et al. A novel coronavirus from patients
300 with pneumonia in china, 2019. New England journal of medicine, 2020.

- 301 [9] Naveena Yanamala, Nanda H Krishna, Quincy A Hathaway, Aditya Radhakrishnan,
302 Srinidhi Sunkara, Heenaben Patel, Peter Farjo, Brijesh Patel, and Partho P Sengupta. A vi-
303 tal sign-based prediction algorithm for differentiating covid-19 versus seasonal influenza
304 in hospitalized patients. NPJ digital medicine, 4(1):1–10, 2021.
- 305 [10] Prasith Baccam, Catherine Beauchemin, Catherine A Macken, Frederick G Hayden, and
306 Alan S Perelson. Kinetics of influenza a virus infection in humans. Journal of virology,
307 80(15):7590–7599, 2006.
- 308 [11] Antonio Gonçalves, Julie Bertrand, Ruian Ke, Emmanuelle Comets, Xavier De Lambal-
309 lerie, Denis Malvy, Andrés Pizzorno, Olivier Terrier, Manuel Rosa Calatrava, France Men-
310 tré, et al. Timing of antiviral treatment initiation is critical to reduce sars-cov-2 viral load.
311 CPT: pharmacometrics & systems pharmacology, 9(9):509–514, 2020.
- 312 [12] Pengxing Cao, Ada WC Yan, Jane M Heffernan, Stephen Petrie, Robert G Moss, Louise A
313 Carolan, Teagan A Guarnaccia, Anne Kelso, Ian G Barr, Jodie McVernon, et al. Innate
314 immunity and the inter-exposure interval determine the dynamics of secondary influenza
315 virus infection and explain observed viral hierarchies. PLoS computational biology,
316 11(8):e1004334, 2015.
- 317 [13] Adrienne L Jenner, Rosemary A Aogo, Sofia Alfonso, Vivienne Crowe, Xiaoyan Deng,
318 Amanda P Smith, Penelope A Morel, Courtney L Davis, Amber M Smith, and Morgan
319 Craig. Covid-19 virtual patient cohort suggests immune mechanisms driving disease out-
320 comes. PLoS pathogens, 17(7):e1009753, 2021.
- 321 [14] Finlay McNab, Katrin Mayer-Barber, Alan Sher, Andreas Wack, and Anne O’garra. Type
322 i interferons in infectious disease. Nature Reviews Immunology, 15(2):87–103, 2015.
- 323 [15] Frederick G Hayden, R Fritz, Monica C Lobo, W Alvord, Warren Strober, Stephen E
324 Straus, et al. Local and systemic cytokine responses during experimental human influenza
325 a virus infection. relation to symptom formation and host defense. The Journal of clinical
326 investigation, 101(3):643–649, 1998.
- 327 [16] Kasia A Pawelek, Giao T Huynh, Michelle Quinlivan, Ann Cullinane, Libin Rong, and
328 Alan S Perelson. Modeling within-host dynamics of influenza virus infection including
329 immune responses. PLoS computational biology, 8(6):e1002588, 2012.
- 330 [17] Naveen K Vaidya, Angelica Bloomquist, and Alan S Perelson. Modeling within-host
331 dynamics of sars-cov-2 infection: A case study in ferrets. Viruses, 13(8):1635, 2021.
- 332 [18] Licia Bordi, Giuseppe Sberna, Eleonora Lalle, Pierluca Piselli, Francesca Colavita,
333 Emanuele Nicastrì, Andrea Antinori, Evangelo Boumis, Nicola Petrosillo, Luisa Mar-
334 chioni, et al. Frequency and duration of sars-cov-2 shedding in oral fluid samples assessed
335 by a modified commercial rapid molecular assay. Viruses, 12(10):1184, 2020.

- 336 [19] Waleed H Mahallawi, Ali Dakhilallah Alsamiri, Alaa Faisal Dabbour, Hamdah Alsaedi,
337 and Abdulmohsen H Al-Zalabani. Association of viral load in sars-cov-2 patients with
338 age and gender. Frontiers in Medicine, 8:39, 2021.
- 339 [20] Keisuke Ejima, Kwang Su Kim, Christina Ludema, Ana I Bento, Shoya Iwanami, Ya-
340 suhisa Fujita, Hirofumi Ohashi, Yoshiki Koizumi, Koichi Watashi, Kazuyuki Aihara,
341 et al. Estimation of the incubation period of covid-19 using viral load data. Epidemics,
342 35:100454, 2021.
- 343 [21] Robert Tibshirani. Regression shrinkage and selection via the lasso. Journal of the Royal
344 Statistical Society: Series B (Methodological), 58(1):267–288, 1996.
- 345 [22] Xiaoyuan Han, Mohammad S Ghaemi, Kazuo Ando, Laura S Peterson, Edward A
346 Ganio, Amy S Tsai, Dyani K Gaudilliere, Ina A Stelzer, Jakob Einhaus, Basile Bertrand,
347 et al. Differential dynamics of the maternal immune system in healthy pregnancy and
348 preeclampsia. Frontiers in immunology, 10:1305, 2019.
- 349 [23] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences.
350 Artificial intelligence, 267:1–38, 2019.
- 351 [24] Been Kim, Rajiv Khanna, and Oluwasanmi O Koyejo. Examples are not enough, learn
352 to criticize! criticism for interpretability. Advances in neural information processing
353 systems, 29, 2016.
- 354 [25] W James Murdoch, Chandan Singh, Karl Kumbier, Reza Abbasi-Asl, and Bin Yu. Defi-
355 nitions, methods, and applications in interpretable machine learning. Proceedings of the
356 National Academy of Sciences, 116(44):22071–22080, 2019.
- 357 [26] Christoph Molnar. Interpretable machine learning. Lulu. com, 2020.