

# Predicting disease transmission rates for hybrid modeling of epidemic outbreaks: statistical and machine learning approaches

Maria Koshkareva<sup>1</sup>[0000–0002–4430–4503], Elizabetty Guseva<sup>1</sup>, Alyona Sharova<sup>1</sup>[0009–0009–4615–7592], and Vasiliy Leonenko<sup>1</sup>[0000–0001–7070–6584]

ITMO University, Kronverksky Pr. 49A, 197101 St. Petersburg, Russia  
{mpkoshkareva,vnleonenko}@itmo.ru

**Abstract.** Hybrid disease modeling is a perspective area of research that allows using detailed individual-based models for the outbreak onset phase and lightweight compartmental models to capture the general trend of the disease progression. In such a way, the method of hybrid modeling provides a good trade-off between the simulation speed and the accuracy of reproducing disease dynamics. One of the problems related to this approach is how to switch properly between the two models. That included detecting the right time moment to finish simulations with the detailed model and calculating correctly the input parameters for the compartmental model. In this paper, we propose an implementation of switching which relies on evaluation and prediction of disease transmission rate. Using an example with a network-based model and a discrete compartmental model, we demonstrate several methods of disease transmission prediction based on statistical models and machine learning approaches and analyze their advantages and disadvantages. The developed methods can be generalized to hybrid modeling of highly detailed demographic processes and propagation processes in general.

**Keywords:** network models · SEIR models · epidemics · compartmental models · machine learning · statistical models · disease outbreaks

## 1 Introduction

For epidemic outbreak modeling, compartmental models, such as SIR [6] and its modifications, are widely used. Their main drawback is a simplifying assumption of homogeneous mixing among individuals. Detailed individual-based models, like agent-based models (ABM) [7] or network-based models [10], can provide more realistic interactions of individuals allowing to capture localized transmission and super-spreading events. However, detailed models have higher computational costs and execution time, which is not desirable for situations requiring fast decision-making. As a result, selecting the best approach between the two for a particular use case is not an easy task [11]. Good news lies in the fact that the advantages of mentioned model techniques can be combined and drawbacks can be compensated by applying hybrid approaches.

Hybrid frameworks, such as those proposed by Bobashev et al. [3] and Hunter et al. [5], have shown that the usage of interacting simple and detailed models can improve both simulation accuracy and computational efficiency. The switching between the submodels occurs in a certain phase of the epidemic. It is assumed that in the early phase of the outbreak, when the number of infected individuals is small, detailed modeling of contacts plays a significant role, which makes an individual-based model a better option for simulation. Particularly, if transmission heterogeneity is high across a pathogen group, a small number of individuals plays a disproportionate role to the spread of a pathogen and targeting control measures towards those individuals can be very effective to reduce the epidemic burden [13]. Once the number of the infected is large, the population can be considered homogeneous, so transitioning to a compartmental model can be done. A key challenge is to choose the optimal switching moment. Early switching leads to information loss related to details of disease transmission, while late switching reduces the computational advantages of the hybrid approach. The success of hybrid simulation depends on an accurate switching moment, as well as on a proper alignment of submodels to ensure that the modeled disease prevalence curve won't show sudden surges or drops in the moment of switching [8].

Among the parameters of epidemic submodels, the disease transmission rate  $\beta$  plays a major role and requires special attention when switching between different model types. As it is noted in [3], thanks to the law of large numbers, the calculated cumulative average of  $\beta$  tends to stabilize as the number of the infected gets larger, which indicates a proper moment to use a compartmental model (a switch point) without altering epidemic dynamics. Hence, one can detect a switch point by waiting for the cumulative average of  $\beta$  to become constant and use this value further on in the simulation via the compartmental submodel. The authors of [3] pay attention to the fact that their hybrid model is built on top of a rather simple ABM (homogeneously mixed population, no community structure) which might lower the effectiveness of their approach for more complicated ABMs. Particularly, other works, such as [12], indicate that assuming time-dependent  $\beta = \beta(t)$  throughout the whole simulation, rather than approximating it with a constant, is essential for accurately recreating epidemic trajectories in real settings.

In this research, we develop and assess methods of switch point detection and  $\beta$  estimation to ensure accurate switching between the submodels. As a baseline method, we regarded the switching of models using a last estimated value of  $\beta$ . The alternative proposed methods include various ways of dynamic  $\beta$  estimation along with  $\beta$  prediction based on the historically known epidemic waves. The methods are tested on a hybrid model which couples a detailed network submodel with a lightweight compartmental submodel. The goal of this work is to compare the computational efficiency (runtime) and accuracy (RMSE) of the analyzed methods. The source code implementing the methods proposed in this paper is available on GitHub [1].

## 2 Methods

### 2.1 Submodels for the hybrid approach

**Compartmental submodel.** The baseline compartmental submodel we use in this study is a deterministic discrete SEIR model with a time step equal to one day. To address the stochasticity which is intrinsic to network and agent-based submodels of the hybrid framework, we also employ the stochastic version of the same model with probabilistic transitions using binomial sampling. In the experiments further on we demonstrate both the deterministic and the stochastic compartmental models. Let  $S$  be susceptible individuals,  $E$  — exposed individuals,  $I$  — infectious individuals, and  $R$  — recovered individuals. The dynamics of the groups’ sizes over time for a discrete stochastic case are set by the following difference equations:

$$\begin{aligned}
 S_{t+1} &= S_t - \xi_{SE}, & (1) \\
 E_{t+1} &= E_t + \xi_{SE} - \xi_{EI}, \\
 I_{t+1} &= I_t + \xi_{EI} - \xi_{IR}, \\
 R_{t+1} &= R_t + \xi_{IR}, \\
 S_0 \geq 0, E_0 \geq 0, I_0 \geq 0, R_0 \geq 0, \\
 S_0 + E_0 + I_0 + R_0 &= N, & (2)
 \end{aligned}$$

where  $\xi_{SE} \sim \text{Bin}(S_t, \beta I_t)$ ,  $\xi_{EI} \sim \text{Bin}(E_t, \sigma)$ ,  $\xi_{IR} \sim \text{Bin}(I_t, \gamma)$ .

**Network submodel.** The population is modeled as a network where each individual is a node, and interactions between the individuals are possible if and only if they are connected by an edge. In our simulations, we used a Barabasi-Albert network topology [2]. The Barabasi-Albert model generates networks through preferential attachment: when new nodes are added, they preferentially connect to existing nodes that already have a high degree (many connections). This process leads to the emergence of “hubs” – highly connected individuals – and a long-tailed degree distribution. In the context of disease spread, this implies that infections may spread more rapidly through these hubs. The network submodel does not have the  $\beta$  parameter as in the SEIR model. Instead, the disease transmission is characterized by  $\tau$ , the probability for one person to infect the other via the common edge.

In the network submodel, the infection dynamics is defined by the contact network topology and the value of transmission probability  $\tau$ . In a compartmental model, the infection transmission is solely governed by the transmission rate  $\beta$ . We assess  $\beta_t$  from a modeled output generated by the network model, using an approximate formula (Eq. 3):

$$\beta_t = -\frac{S_{t+1} - S_t}{S_t \cdot I_t}. \quad (3)$$

As stated earlier, the choice of  $\beta$  values for the compartmental submodel is important to ensure a proper switch. To address this, we generated baseline

epidemic prevalence trajectories made solely by the network submodel, without switching. The values of  $\sigma$  and  $\gamma$  were fixed, we varied only 2 parameters –  $\tau$  and  $I_0$ , the fraction of initially infected. Resulted epidemic curves, as will be described in detail later, constituted test and train datasets which were used to evaluate our methods of  $\beta$  analysis and prediction.

## 2.2 Beta estimation approaches

The regarded approaches were grouped into three categories: (a) those that rely on a current incomplete disease trajectory (which mimics the ongoing epidemic process), (b) those that are based on a train dataset with complete trajectories (which helps establish the form of  $\beta_t$  based on previously collected information), and (c) those that use both a train set and a current disease trajectory (which makes it possible to better adjust a form of  $\beta_t$  to an actual disease curve). Each group of approaches has their own baselines. In graphs and tables, we use the abbreviated names in the form of  $M_{\langle type \rangle}$ , where  $\langle type \rangle$  is the method name.

**Estimation on current incomplete data.** To estimate the value of  $\beta$  for switching, we started with three baselines: choosing the last known  $\beta_t$  value, i.e. the closest to a switch point ( $M_{last\_val}$ ); taking the last value from the moving averages of  $\beta_t$  ( $M_{ma\_val}$ ), taking the last value from the cumulative averages of  $\beta_t$  ( $M_{ca\_val}$ ). The more advanced method of  $\beta$  estimation is to fit a function to incomplete  $\beta_t$  data before the day of switch and use it in a compartmental model ( $M_{biexp}$ ). The function should have a similar shape with actually observed trends of  $\beta_t$ , i.e. bear similarity with a skewed bell. The chosen function is the biexponential decay function (Eq. 4):

$$\beta_t = a(e^{-bt} - e^{-ct}), \quad (4)$$

where  $a$ ,  $b$  and  $c$  are estimated through non-linear least squares.

**Estimation on train set.** The methods in this group obtain the form of  $\beta_t$  based on the train set, therefore the form of  $\beta_t$  is fixed and is not influenced by actual prevalence data related to the ongoing simulated outbreak. The baseline method ( $M_{median}$ ) for this group is to use the trajectory of median values of  $\beta_t$  for each day  $t$  ( $\hat{\beta}_t$ ) from the train set. The alternative method is to use  $\beta_t$  in the form of third order polynomial regression ( $M_{regr}$ ) with L2 regularization. The model takes  $t$  as input and outputs  $\log(\beta)$ . Input values are modified by removing the mean and scaling to variance.

**Estimation on train set and incomplete data.** The methods from this group use both generated trajectories and the current sample’s data to generate the forecast for the model switch. The baseline methods consist of taking the forecast  $\beta$  values from the previous group of methods and making them better comply with  $\beta$  values of current data.

The first baseline method ( $M_{median}^{shift}$ ) is to shift modeled values of  $M_{median}$ , adding a scalar. The summand is calculated as the difference between an actual and estimated  $\beta$  at the switch point. The value of the actual  $\beta$  is chosen as a moving average at the switch point. Let  $\beta'$  denote  $\beta$  estimate from the baseline and  $t$  – the time of switch, then the shifted trajectory is set by the following equation:

$$\beta'_{shifted} = \beta' + (MA_t - \beta'_t), \quad (5)$$

where  $MA_t$  is a moving average for  $\beta_t$ .

The second baseline ( $M_{regr}^{shift}$ ) is to shift a regression forecast ( $M_{regr}$ ) in a similar manner. The third one ( $M_{regr}^{add}$ ) is to additionally train a regression model for additional  $N$  epochs, using  $\beta$  values from the currently observed outbreak as target values.

The main methods are a regression model with an extended set of income features ( $M_{regr\_ext}$ ) and a Long Short-Term Memory (LSTM) network ( $M_{LSTM}$ ) [4]. LSTM was chosen as it is often used for time series forecasting and for epidemic modeling in particular [9].

The regression model is with third degree polynomial features, i.e. all polynomial combinations of the features with degree 1 and 2. The regularization is L2. Input values are scaled by removing the mean and scaling to unit variance.

The LSTM network consists of 2 LSTM layers. The loss function is a mean squared error (MSE), the optimizer is RMSprop with an initial learning rate of 0.001 and a learning rate schedule: multiply the learning rate by 0.1 after 30 epochs. The training was conducted with a batch size of 64 for 100 epochs with early stopping: training was stopped if validation loss did not improve after 7 epochs. The model was trained to forecast based on 14 days prior. The chosen architecture (Fig. 1) is complex enough to catch the trends for  $\beta$  prediction on the train dataset.

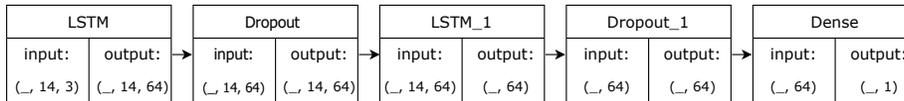


Fig. 1. LSTM architecture

### 2.3 Switch point detection approaches

To determine the optimal time of switch from the network submodel to the compartmental submodel, we need to monitor  $\beta$  values. Initial days of an epidemic are characterized by a highly variable  $\beta$  with its further smoothing, therefore our criteria for switching is low  $\beta$  variance. This ensures that the compartmental model can accurately reproduce epidemic dynamics without significant loss of detail.

To compare dynamic methods described further, we also utilize static methods with constant switch points, later referred to in the form of  $SP_{<n>}$ , where  $n$  is the chosen day of switch.

**Change in variance.** The first method ( $SP_{once}$ ) involves calculating the moving variance of  $\beta$ , i.e. calculating variance over the previous  $k$  days. Let  $MV_t$  denote the calculated moving variance value for a given day  $t$ :

$$MV_t = Var(\beta_{t-k+1}, \dots, \beta_t) \quad (6)$$

To establish a switch condition regardless of the range of  $\beta$  values, we apply min-max scaling to get values in the range from 0 to 1 (Eq. 7):

$$MV'_t = \frac{MV_t - MV_{min}}{MV_{max} - MV_{min}} \quad (7)$$

This scaling allows us to interpret  $MV'_t = 0$  as the minimal observed variability. Therefore, when  $MV'_t$  falls below a chosen threshold ( $\varepsilon$ ), we can consider  $\beta$  stable, which indicates the optimal time to switch.

**Established change in variance.** At the beginning of an epidemic the variance of  $\beta$  may accidentally decrease, not reflecting real stabilization. To account for such declines leading to early switching, we add the following condition: the value of  $MV'_t$  should stay below  $\varepsilon$  for a specified number of consecutive days. Therefore the second method ( $SP_{estab}$ ) waits for an established change in variance.

**Change in variance after the epidemic situation.** The third method ( $SP_{epi\_sit}$ ) introduces a new condition to avoid the early switch problem. The first condition, i.e. making sure that  $MV'_t$  is below  $\varepsilon$ , is accompanied by a second condition: the proportion of the infected population should reach a specified percent; this indicates the epidemic situation.

## 2.4 Evaluation metrics

To assess the computational efficiency of the algorithms, time spent from the model initialization to the final prediction was measured. The observed period includes both training and inference phases. Forecast accuracy was measured by calculating RMSE between the actual and the predicted time series, both for  $I_t$  and for  $\beta_t$ .

Two additional metrics for  $I_t$  are peak height error (relative) and peak time error (absolute). Peak time error is negative when the model peak day is earlier than the real one, and positive otherwise. The value of peak height corresponds to maximal burden imposed on the healthcare units, whereas the value of peak time is used to assess the period of time remaining to prepare necessary resources. In case of positive peak time error, we will be caught by surprise unprepared, that is why negative peak time error is preferable.

### 3 Experiments and results

In our simulations, we used a Barabasi-Albert graph network of  $10^4$  nodes, with each new node attached to 5 existing nodes. Epidemic trajectories were generated with  $\sigma = 0.1$ ,  $\gamma = 0.08$ ,  $\tau \in \{0.01, 0.02, \dots, 0.05\}$ ,  $I_0 \in \{0.005, 0.006, \dots, 0.011\}$ . We regard these trajectories as outbreaks of some generic acute respiratory disease. The values of parameters  $\sigma$  and  $\gamma$  were chosen to be close to typical incubation and infection periods for ARIs. Values of  $\tau$  and  $I_0$  were selected to approximately match ARI outbreaks by the length of an epidemic and the maximal prevalence.

For each parameter combination, 50 simulation runs were performed with different seeds of a random number generator, thus reflecting the stochastic factor. A total of 1 500 unique epidemic trajectories were generated, with 20% of the data selected as a test dataset using a stratified approach based on the parameter values  $\tau$  and  $I_0$ . All further experiments test the methods of  $\beta$  estimation and switch point detection on these 300 samples. The results for RMSE of  $I_t$  are presented in Table 1 and discussed further.

$\beta$ estimation method	Switch point detection method					
	$SP_{20}$	$SP_{30}$	$SP_{40}$	$SP_{once}$	$SP_{estab}$	$SP_{epi\_sit}$
$M_{last\_val}$	244.8 ± 116.0	154.7 ± 134.8	99.7 ± 134.8	265.2 ± 111.0	265.1 ± 109.1	287.5 ± 98.9
$M_{ma\_val}$	273.7 ± 88.5	199.2 ± 124.6	111.5 ± 132.6	264.4 ± 84.6	262.4 ± 83.0	291.7 ± 86.9
$M_{ca\_val}$	229.9 ± 59.9	221.0 ± 74.7	193.3 ± 95.3	244.8 ± 68.5	248.3 ± 69.7	243.3 ± 66.3
$M_{bixep}$	261.4 ± 152.8	205.9 ± 121.2	153.8 ± 84.8	324.8 ± 206.1	290.6 ± 202.4	325.1 ± 190.1
$M_{median}$	<b>74.2 ± 90.1</b>	<b>42.3 ± 56.9</b>	<b>28.6 ± 38.8</b>	<b>87.2 ± 76.2</b>	<b>83.8 ± 66.0</b>	<b>85.6 ± 71.6</b>
$M_{regr}$	109.6 ± 75.8	50.5 ± 65.5	42.7 ± 49.0	139.3 ± 77.4	131.5 ± 74.0	144.0 ± 73.5
$M_{median}^{shift}$	162.9 ± 117.0	70.9 ± 74.3	34.6 ± 61.5	178.1 ± 156.8	161.2 ± 159.6	183.6 ± 136.2
$M_{regr}^{shift}$	213.5 ± 67.0	112.5 ± 62.8	45.2 ± 71.6	209.1 ± 86.0	196.8 ± 87.9	217.2 ± 77.5
$M_{regr}^{add}$	83.4 ± 53.4	59.8 ± 48.9	48.2 ± 42.8	125.0 ± 61.1	120.0 ± 59.6	128.0 ± 58.1
$M_{regr\_ext}$	<b>58.1 ± 48.2</b>	<b>37.2 ± 36.8</b>	<b>28.2 ± 28.5</b>	<b>66.9 ± 62.0</b>	<b>64.0 ± 60.2</b>	<b>67.1 ± 56.5</b>
$M_{LSTM}$	<b>80.9 ± 68.5</b>	<b>31.6 ± 26.8</b>	<b>20.9 ± 13.6</b>	<b>102.7 ± 85.6</b>	<b>94.6 ± 86.3</b>	<b>105.1 ± 85.5</b>

**Table 1.** RMSE of  $I_t$  for all methods of beta estimation and switch point detection; top-3 beta estimation methods for each switch point are highlighted in bold

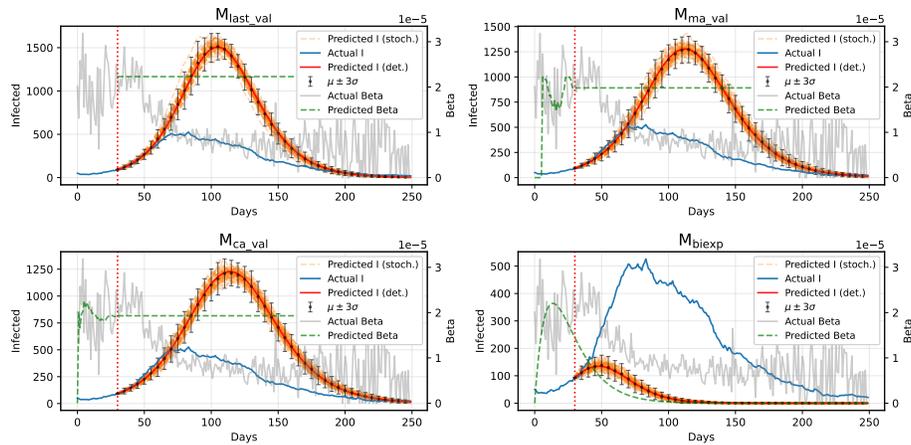
#### 3.1 Beta estimation approaches

The aim of these experiment series is to compare the accuracy of  $\beta$  estimation approaches within each group, corresponding to their use cases.

**Estimation on current incomplete data.** We applied the methods of  $\beta$  estimation from the first group (section 2.2) to each test sample. Our additional subject of interest is to compare prediction results for constant and time-dependent  $\beta$ .

As can be seen in Fig. 2 and first 3 rows in Table 1, using constant  $\beta$  as input for the SEIR provides a poor fit to actual values of  $I_t$ . The  $\beta$  trajectory is highly variable, so using  $M_{last\_val}$  leads to a big difference in RMSE based on the day of switch. The moving average ( $M_{ma\_val}$ , the best results with a window size of 7 days) and the cumulative average ( $M_{ca\_val}$ ) approaches result in lower variance of RMSE compared to  $M_{last\_val}$ . To further analyze limitations of constant  $\beta$  values, we took values from 0 to  $4 \cdot 10^{-5}$  with a step size  $10^{-6}$  as inputs to the SEIR model. The modeled  $I_t$  trajectories either are too wide or have a higher peak to match actual  $I_t$  (Fig. 3, left). This suggests that constant  $\beta$  gives poor results because the data requires a varying  $\beta$  trajectory, not because the constant value was chosen incorrectly.

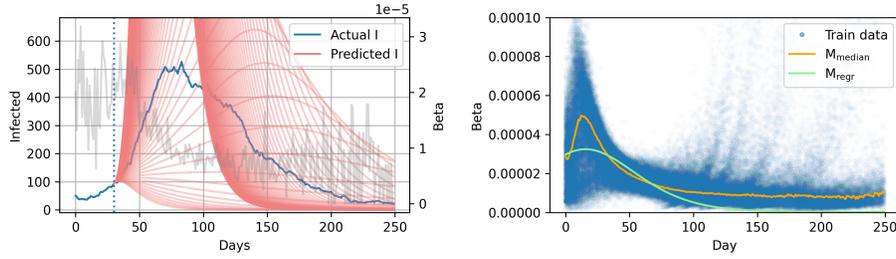
The first time-dependent  $\beta$  estimation approach ( $M_{biexp}$ ) does not estimate values similar to the initial  $\beta$ . The modeled  $\beta$  trajectory after the switch point does not match the curve of real values either, which leads to peak height underestimation. However, the method with time-dependent  $\beta$  has area of improvement, which will be showed further.



**Fig. 2.** Simulated prevalence curves for beta estimation methods based on current data: last value, moving average, cumulative average, biexponential decay

**Estimation on train set.** We applied the methods of  $\beta$  estimation from the second group (section 2.2) to each test sample.

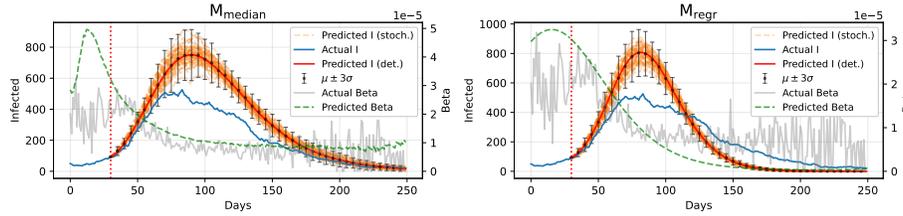
As can be seen in Fig. 3 (right), the maximum variation of actual values occurs at the beginning of an epidemic. Then, when an outbreak starts spreading rapidly and  $\beta$  values have the most effect on modeled  $I_t$  trajectories, the range of values is narrower. The methods  $M_{median}$  and  $M_{regr}$  have similar estimated trajectories after the 30<sup>th</sup> day when  $\beta$  values should be the most relevant for modeling, but the metrics vary significantly with  $M_{median}$  having lower errors.



**Fig. 3.** Modeled trajectories for  $I_t$  with constant beta values (left); median  $\beta$  values of all generated trajectories and the output of the regression model fit on a day (right)

This implies that a slight mismatch even solely in first initial  $\beta$  values can nearly double the median RMSE of  $I_t$ .

Even a slight difference in  $\beta$  values results in a noticeable change in a number of infected people. This can be supported by Fig. 3 (left): each increment by  $10^{-6}$  gives around 100 additional infected people at the modeled peak, i.e. 1% of the population.



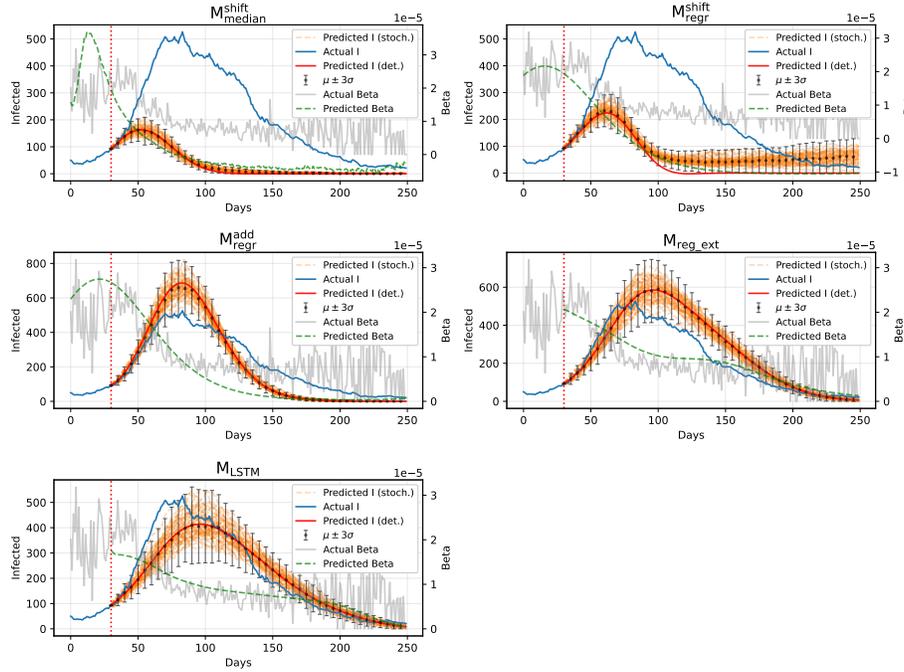
**Fig. 4.** Simulated prevalence curves for beta estimation methods based on train set: median of train set, regression on day

**Estimation on train set and incomplete data.** We applied the methods of  $\beta$  estimation from the third group (section 2.2) to each test sample.

Both shift methods (the best results with a window size of 14 days) result in worse metrics than original methods ( $M_{median}$  and  $M_{regr}$ ). One reason may be the shape of  $\beta$  values in simulated trajectories. While initial days have a high variability in values, shifting based on these days' values (even with rolling average) results in inflated  $\beta$  values at the end.

We can compare  $M_{regr}^{add}$  (the best results with 3 additional epochs) in Fig. 5 with  $M_{regr}$  in Fig. 4. The tail of estimated  $\beta$  from  $M_{regr}^{add}$  is the same, but the beta values near the switch point are closer to the actual  $\beta$ . This gives a lower peak, although the epidemic generally has the same duration.

$M_{reg\_ext}$  and  $M_{LSTM}$  methods are in top-3 for RMSE of  $I_t$  according to Table 1. Input features for  $M_{reg\_ext}$  and  $M_{LSTM}$  were selected to achieve the least RMSE without redundant features. For  $M_{reg\_ext}$ , the input features are:  $t, S_t, E_t, I_t, R_t, I_{t-1}$ . The input features for  $M_{LSTM}$  are:  $t, E_t, I_{t-2}$ . Feature values at switch point are compartment values from the network model. The next day is modeled with compartments from the discrete SEIR. We experimented with different values for  $I_{t-s}$ , where  $s$  is the shift;  $s = 2$  gave the best results.



**Fig. 5.** Simulated prevalence curves for beta estimation methods based on current data and train set: shifted median, shifted regression, regression with additional learning, regression with extended input features, LSTM

### 3.2 Switch point detection approaches

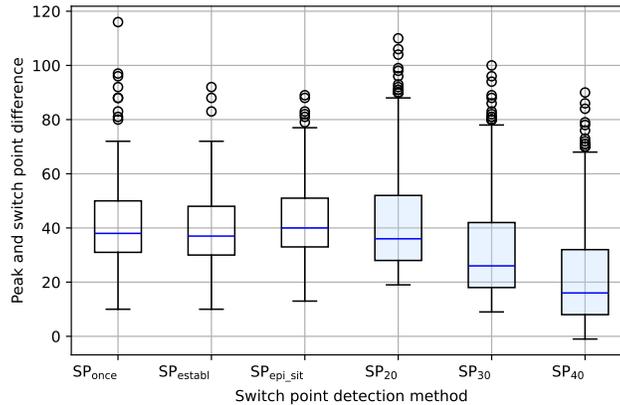
The aim of these experiment series is to assess switch point detection approaches and the stability of  $\beta$  estimation methods. To achieve this, we calculated the metrics for a fixed set of switch points and for variance-based detection approaches. RMSE of  $I_t$  for all  $\beta$  estimation methods and switch point detection methods can be found in Table 1.

The lowest accepted switch point was set to 14 days to not interfere with  $M_{LSTM}$ 's input shape. The most test samples visually have an epidemic start

(1% of the population is infected) around days 20, 30, and 40. We chose these three days as constant switch points. The choice of switch point does not affect the top-3  $\beta$  estimation methods based on RMSE. This may suggest that even though  $SP_{20}$  may be too early to switch, or  $SP_{40}$  is too close to the peak, no method of  $\beta$  estimation gains any benefits to overcome the top-3 methods ( $M_{median}$ ,  $M_{regr\_ext}$  and  $M_{LSTM}$ ).

Switch detection methods have a different optimal threshold  $\varepsilon$ .  $SP_{once}$  shows best results with  $\varepsilon = 0.05$ , i.e. 5%.  $SP_{estab}$  has best results with  $\varepsilon = 0.05$  and 2 consecutive days.  $SP_{epi\_sit}$  shows best results with  $\varepsilon = 0.1$ . The discrepancy comes from the methods' conditions: the last approach waits for 1% of the population to be infected, further  $\beta$  values are more stable, so we can switch with fewer risks, therefore a higher  $\varepsilon$ .

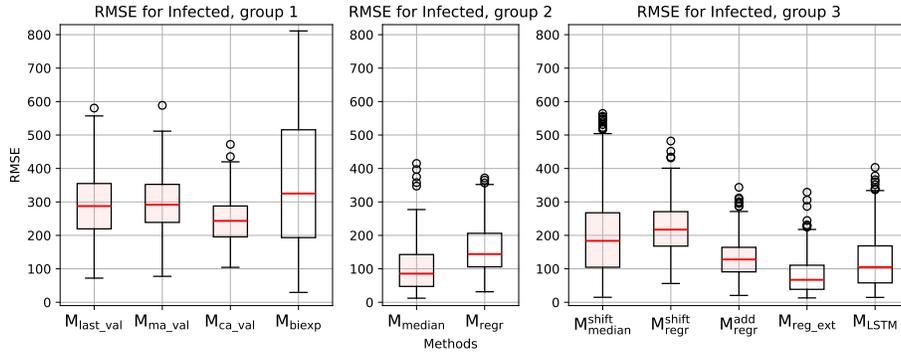
Constant switch points have limited usage for real-world applications because they do not consider the ongoing epidemic dynamics. Out of the three remaining methods,  $SP_{estab}$  has the lowest errors, according to Table 1. However, high accuracy may be a result of switching closer to the peak, which reduces the benefits of the hybrid approach. For instance, as it can be seen in Table 1, increasing the value of a constant switch point (from 20 to 30 and 40) results in decreasing the RMSE of  $I_t$  for the corresponding methods ( $SP_{20}$ ,  $SP_{30}$ ,  $SP_{40}$ ). Thus, it is important to also consider the distance between the detected switch point and the actual peak time, i.e. the difference between their values in days. Fig. 6 shows that the method with the largest median difference between the switch and the peak time is  $SP_{epi\_sit}$ . Due to that reason, we prefer this method, although it has the largest errors among the 3 methods based on the variance of  $\beta$ .



**Fig. 6.** Difference between the actual peak time and the detected switch point; boxplots for constant switch points are filled

### 3.3 Evaluation results

**Accuracy (RMSE).** All figures are presented for the best switch point detection approach –  $SP_{epi\_sit}$ . For RMSE of  $I_t$  (Fig. 7), the methods with the best results are:  $M_{LSTM}$ ,  $M_{median}$  and  $M_{regr\_ext}$  with median errors 105.1, 85.5 and 67.1. Therefore, the error for top-3 methods is around 1% of the population. For RMSE of  $\beta$  values (Fig. 8), the best performing methods are:  $M_{LSTM}$ ,  $M_{median}$  and  $M_{regr\_ext}$  with median errors  $1.09 \cdot 10^{-5}$ ,  $1.08 \cdot 10^{-5}$  and  $1.05 \cdot 10^{-5}$ . The set of best methods is the same for both metrics. The purpose of a hybrid approach is to switch to a simpler model without the loss in accuracy for the predicted  $I_t$  trajectory, so RMSE of  $\beta$  is less relevant. As a result, the final top-3 methods are:  $M_{LSTM}$ ,  $M_{median}$ ,  $M_{regr\_ext}$ .

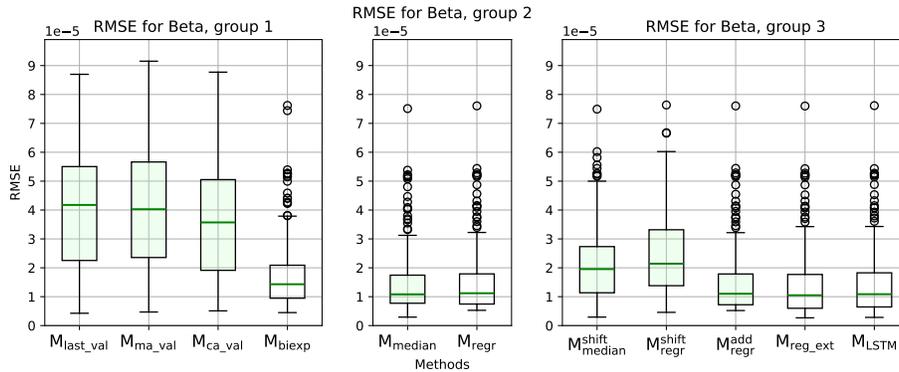


**Fig. 7.** RMSE for  $I_t$ . Beta estimation methods based on current data (group 1), train set (group 2), and combined (group 3). Boxplots for baseline methods are filled.

	$M_{last\_val}$	$M_{ma\_val}$	$M_{ca\_val}$	$M_{biexp}$	$M_{median}$	$M_{regr}$
Time, sec	$0.016 \pm 0.006$	$0.005 \pm 0.005$	$0.003 \pm 0.005$	$0.003 \pm 0.006$	$0.008 \pm 0.006$	$0.005 \pm 0.001$
	$M_{shift\_median}$	$M_{shift\_regr}$	$M_{add\_regr}$	$M_{regr\_ext}$	$M_{LSTM}$	
Time, sec	$0.009 \pm 0.007$	$0.005 \pm 0.000$	$0.007 \pm 0.001$	$0.362 \pm 0.073$	$37.625 \pm 14.519$	

**Table 2.** Prediction time for the methods of  $\beta$  estimation

**Computational efficiency.** Time measurements for training and prediction were performed on a system equipped with an AMD Ryzen 5 5500U processor (up to 4.0 GHz), 16 GB of LPDDR4x RAM (4266 MHz), and integrated AMD Radeon Graphics. In a single-threaded configuration without parallelization,  $M_{LSTM}$  required 1228.4 seconds to finish the training,  $M_{regr\_ext}$  took 2.86



**Fig. 8.** RMSE for Beta. Beta estimation methods based on current data (group 1), train set (group 2), and combined (group 3). Boxplots for baseline methods are filled.

seconds. Time for predictions is presented in Table 2.  $M_{LSTM}$ , while demanding the most computational resources, is still considered top-3 due to its forecast accuracy. Two methods with higher RMSE and faster prediction time are  $M_{median}$  and  $M_{regr\_ext}$ .

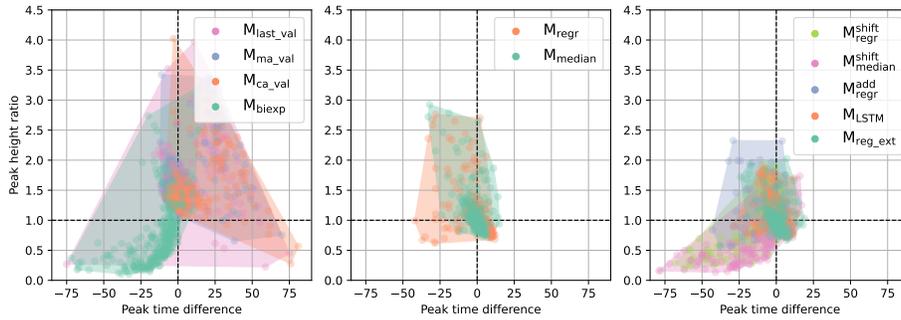
**Accuracy (distance to peaks).** All figures are presented for the best switch point detection approach –  $SP_{epi\_sit}$ .

For an easier interpretation, we present peak errors for each group of  $\beta$  estimation approach separately (Fig. 9). The x-axis corresponds to the difference between the predicted and actual peak time. The y-axis is the fraction of the predicted peak height to the actual peak height. The best case is at the point with coordinates (0, 1). Besides the best case, we are also interested in the second and third quadrants, which depict cases with peaks predicted earlier. All top-3  $\beta$  estimation methods ( $M_{median}$ ,  $M_{LSTM}$  and  $M_{regr\_ext}$ ) are close to the best case at (0, 1).

The remaining  $\beta$  estimation methods have larger areas of possible outcomes. Some additional conclusions may be drawn from methods based on regression ( $M_{regr}$ ,  $M_{regr}^{shift}$ ,  $M_{regr}^{add}$ ). Shifted predictions by  $M_{regr}^{shift}$  skew the results closer to the best case for peak height ratio but with larger peak time underestimation. Additional learning for 3 epochs in  $M_{regr}^{add}$  makes the error in peak height slightly lower for almost each point.

## 4 Conclusions and future work

In this research, we presented the hybrid approach based on a network-based and a discrete SEIR submodels with dynamic switching. For the correct alignment of submodels, it is important to properly estimate the value of  $\beta$  (disease transmission rate). We conducted several experiments to analyze methods of  $\beta$



**Fig. 9.** Peak errors for beta estimation methods; x-axis is the difference between the predicted and actual peak time, y-axis is the fraction of the predicted peak height to the actual peak height

estimation and switch point detection.  $\beta$  estimation methods are divided into 3 groups based on their use conditions: only the current outbreak data are available, only historic data are available, historic and current data are available. Switch point detection methods are divided into static methods with constant switch points and dynamic detection methods based on  $\beta$  variance. Our work is based on the analysis of synthetic epidemic data; this is the first step towards the further application of the hybrid approach for modeling real-world epidemic processes.

The best switch point detection approach was concluded to be  $SP_{epi\_sit}$ , i.e. switch point detection based on the declared epidemic. The method has two conditions: 1% of  $I_t$  and  $\beta$  variance lower than 10%. The best  $\beta$  estimation methods with the lowest RMSE in  $I_t$  and with the lowest peak errors are:  $M_{LSTM}$  (LSTM model),  $M_{median}$  (median of  $\beta$  values in train set) and  $M_{reg\_ext}$  (regression model).

For future work, firstly, we plan to use interval estimates as opposed to point estimates to account for uncertainty. As some papers suggest, epidemic forecasts should always be done with deep uncertainty methods to enhance decision-making. Secondly, there is a separate field devoted to changepoint detection, which analyses when the change happens in the probability distribution of a certain signal. We can utilize state-of-the-art approaches and assess their applicability for our purposes. Thirdly, to formalize the assessment of switch point detection methods based on accuracy (distance to peak and RMSE of  $I_t$ ) one can include the weighted metric. Finally, we also plan to generalize the methods of  $\beta$  estimation for more complex models such as ABM and real data. However, initial attempts showed more intricate  $\beta$  trajectories for ABM data, thus requiring changes in our approaches. ABM may also have an additional parameter, a fraction of immune individuals, to consider during modeling and  $\beta$  estimation.

**Acknowledgement** This research was supported by The Russian Science Foundation, Agreement #22-71-10067.

## References

1. [https://github.com/vnleonenko/Network\\_hybrid](https://github.com/vnleonenko/Network_hybrid)
2. Albert, R., Barabási, A.L.: Statistical mechanics of complex networks. *Reviews of modern physics* **74**(1), 47 (2002)
3. Bobashev, G.V., Goedecke, D.M., Yu, F., Epstein, J.M.: A hybrid epidemic model: Combining the advantages of agent-based and equation-based approaches. In: 2007 Winter Simulation Conference. pp. 1532–1537 (2007). <https://doi.org/10.1109/WSC.2007.4419767>
4. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Computation* **9**, 1735–1780 (11 1997). <https://doi.org/10.1162/neco.1997.9.8.1735>
5. Hunter, E., Namee, B.M., Kelleher, J.D.: A hybrid agent-based and equation based model for the spread of infectious diseases (2020)
6. Kermack, W.O., McKendrick, A.G.: A contribution to the mathematical theory of epidemics. *Proceedings of the Royal Society of London. Series A, Containing papers of a mathematical and physical character* **115**(772), 700–721 (1927)
7. Kerr, C.C., Stuart, R.M., Mistry, D., Abeysuriya, R.G., Rosenfeld, K., Hart, G.R., Núñez, R.C., Cohen, J.A., Selvaraj, P., Hagedorn, B., et al.: Covasim: an agent-based model of COVID-19 dynamics and interventions. *PLOS Computational Biology* **17**(7), e1009149 (2021)
8. Leonenko, V.: A hybrid modeling framework for city-scale dynamics of multi-strain influenza epidemics. In: *International Conference on Computational Science*. pp. 164–177. Springer (2022)
9. Leonenko, V.N., Bochenina, K.O., Kesarev, S.A.: Influenza peaks forecasting in Russia: Assessing the applicability of statistical methods. *Procedia Computer Science* **108**, 2363–2367 (2017)
10. Pastor-Satorras, R., Vespignani, A.: Epidemic spreading in scale-free networks. *Physical review letters* **86**(14), 3200 (2001)
11. Rahmandad, H., Sterman, J.: Heterogeneity and network structure in the dynamics of diffusion: Comparing agent-based and differential equation models. *Management science* **54**(5), 998–1014 (2008)
12. Smirnova, A., deCamp, L., Chowell, G.: Forecasting epidemics through nonparametric estimation of time-dependent transmission rates using the seir model. *Bulletin of mathematical biology* **81**(11), 4343–4365 (2019)
13. Tran-Kiem, C., Bedford, T.: Estimating the reproduction number and transmission heterogeneity from the size distribution of clusters of identical pathogen sequences. *Proceedings of the National Academy of Sciences* **121**(15), e2305299121 (2024)