# Prediction of the disease controllability in a complex network using machine learning algorithms

Richa Tripathi
  Indian Institute of Technology Gandhinagar

Amit Reza ( ✉ amit.reza@iitgn.ac.in )
  Indian Institute of Technology Gandhinagar    https://orcid.org/0000-0001-7934-0259

Dinesh Garg
  IBM-IRL: IBM India Research Laboratory

# Prediction of the disease controllability in a complex network using machine learning algorithms

Richa Tripathi, Amit Reza and Dinesh Garg

**Index Terms**

Complex Networks, SIR models, Basic reproduction number, Machine Learning.

——————————— ✦ ———————————

## CORRESPONDING AUTHOR

Amit Reza, email id: amit.reza@iitgn.ac.in

## COMPLIANCE WITH ETHICAL STANDARDS

1) Disclosure of potential conflicts of interest: We would like to declare that we have no conflicts of interest.
2) Research involving human participants and/or animals: NA
3) Informed Consent: NA
4) Funding: NA

———————————————————

- *Richa Tripathi, Indian Institute of Technology Gandhinagar, Gandhinagar-382355, Gujarat, India. E-mail: richa.tripathi@iitgn.ac.in*

- *Amit Reza, Indian Institute of Technology Gandhinagar, Gandhinagar-382355, Gujarat, India.*

- *Dinesh Garg, India Research Laboratory, Bangalore, India.*

# Prediction of the disease controllability in a complex network using machine learning algorithms

**Abstract**—The application of machine learning (ML) techniques span a vast spectrum ranging from speech, face and character recognition, medical diagnosis, anomaly detection in data to the general classification, prediction and regression problems. In the present work, we solve the problem of predicting $R_0$ for disease spreading on complex networks using the regression-based state-of-art ML techniques. $R_0$ is a metric that determines whether the disease-free epidemic or an endemic state is asymptotically stable and hence indicates the controllability of the disease spread. We predict $R_0$, based on training the ML models with structural properties of complex networks, irrespective of the network type. The prediction is possible because: (a) The structure of complex networks plays an essential role in the spreading processes on networks (b) The regression techniques such as Support Vector Regression and Artificial Neural Network Model can be very efficiently used for prediction problems, even for non-linear data. We obtained good accuracy in the prediction of $R_0$ for the simulated networks as well as real-world networks using these techniques. Moreover, the ML model training is a one-time investment cost in terms of training time and memory, and the trained model can be used for predicting $R_0$ on unseen/new examples of networks.

**Index Terms**—Complex Networks, SIR models, Basic reproduction number, Machine Learning.

◆

## 1 INTRODUCTION

THe problem of disease spreading has been studied using a system of Ordinary Differential Equations (ODE) [Anderson and May(1992)], to predict the endemic disease state and devise effective control strategies. Earlier studies did not employ any spatial structure, and the dynamics generally depended on the population number and the probabilities of transitions from one disease state to others. However, the use of a spatial structure for determination of relative positions and interactions of individuals has taken a front stage recently. The complex network framework [Barabási and Albert(1999)], where nodes represent the individuals, and the links govern the interactions between the nodes is thus very useful. Disease spreading on networks has been studied using various compartmental epidemiology models such as SI (Susceptible-Infected), SIR (Susceptible-Infected-Recovered), SIRS (Susceptible-Infected-Recovered-Susceptible), etc. [Hethcote(2000)]. The impact of disease in the population is measured using basic reproduction number, $R_0$ [Lotka(1956)]. $R_0$ is the average number of individuals an infected person infects over its period of activity, such that if $R_0 < 1$ the disease will die out in the long run and if $R_0 > 1$ the disease-free stationary state is asymptotically unstable [Stewart et al.(2005)Stewart, Logsdon, and Kelley]. The fact that for a $100\%$ effective vaccine, the fraction of individuals that need to be vaccinated is $1 - \frac{1}{R_0}$ to prevent persistent disease spread, indicates that higher number of individuals need to be vaccinated if the factor $R_0$ is high for a disease. There have been several works [Shirley and Rushton(2005)], [Schimit and Pereira(2018)] for determining the dynamical relationship between network and disease parameters for an epidemic spread occurring on networks.

Machine learning (ML) models based on supervised and unsupervised learning algorithms have found important applications in the area of complex networks. For example, for optimal graph partitioning into community structure [MacQueen et al.(1967)], for classification of diseased networks from the control networks using data from brain imaging [Chaplot et al.(2006)Chaplot, Patnaik, and Jagannathan], for classification of networks into various model networks [Xin et al.(2018)Xin, Zhang, and Shao], etc. Recently, a study [Schimit and Pereira(2018)] reported the relative relevance of network topological and disease parameters for disease spreading on complex networks, irrespective of the model network type. They found that the topology of population (or network) affects the disease spreading process, apart from the disease parameters. In essence, for the given initial conditions for the epidemic spread, the topological properties of networks govern the asymptotic disease state. Our work is based on this idea and focuses on explor-

ing if the topological properties solely could predict $R_0$. The accurate determination/prediction of $R_0$ is of paramount importance to analyze the stability of the disease stage, concerning the infection outbreak. To this end, we train ML regression models, using large number of networks of various model network types. We used six structural properties of five different complex networks examples as input features and the corresponding $R_0$ evaluated after simulating the SIR dynamics on them, as output labels. While training, the model fits these inputs with the outputs and learns certain weights. Using the trained models (or the weights), we predict the $R_0$ value for test network example based on its own structural metrics.

In particular, we trained three models: linear regression, support vector regression (non-linear regression) with different kernels and a neural network model. We optimized these models for the correct prediction of $R_0$ on test examples and evaluated two accuracy metrics: mean squared error and coefficient of determination [Pedregosa et al.(2011)Pedregosa, Varoquaux, Gramfort, Michel, Thirion, Grisel, Blondel, Prettenhofer, Weiss, Dubourg, Vanderplas, Passos, Cournapeau, Brucher, Perrot, and Duchesnay] for each of them. We present the parameter list and their ranges for fine tuning of each of these models. Support Vector Regression (SVR) with radial basis function (RBF) kernel and the artificial neural network (ANN) model resulted in high accuracy of prediction over linear models, for both real and simulated networks. Moreover, the excellent overlap between the expected and predicted values of $R_0$ for these nonlinear ML models also point to the non-linear relationship between the model input and output variables. Hence, we show that simple ML techniques can be used to predict $R_0$ with high precision using the structural properties of the network. We also find that different network metrics have different relative powers of prediction. Based on this result, we can just use four most important measures out of the six, for the model training and prediction. However, ML model performances were marginally better with all the six features used.

## 2 PROBLEM STATEMENT

As explained earlier, the $R_0$ of disease spreading is an estimate of disease impact on the population and hence its controllability. For a dynamical process occurring on a network, the structure of the network plays a key role in determining the next stage of the process apart from inherent parameters of the process. Similarly for a disease spread on a network, the dynamics and hence the disease stages are a result of interplay between the network structural features and the epidemic model parameters such as the probability of infection, probability of recovery from infection,

etc. For the $R_0$ calculation in the present work, we update all the disease related rate constants at each time step depending on the instantaneous values and the respective changes in the populations of the susceptible, infected and recovered individuals according to equations 2-5 in supplementray information (SI). Since in networks the interactions can only occur through nodes neighbours, the instantaneous values and changes of populations of S, I and R over the whole network are indirectly determined by the local and global structure of the network. Hence for networks, the value of $R_0$, which is a dynamical descriptor of the disease and depends on the disease related rate constants is affected by the network structure.

Given that $R_0$ is the average number of susceptible an infected individual infects over its period of being in infected state, different structural metrics of networks have specific effect on its value. For example, for an $Erdos - Rényi$ (ER) network, higher the clustering coefficient means higher the number of connections and hence higher the $R_0$. Similar trend is observed for network density, average degree and the maximum degree. On the other hand, higher the shortest path length means higher is the average shortest distance between two nodes and hence smaller the $R_0$ and vice versa. For a Small-World (SW) network, owing to the regularity of its structure, even the lower density of connections than an ER network shows similar potential for the disease spread. Hence, other topological features are also important for determining the $R_0$. In the same manner, the effect of topological parameters on the $R_0$ value can be intuitively understood for other model networks.

In this work, we seek to predict the $R_0$ for any example network, given that we know its structural properties a priori. Given that $R_0$ is an important measure to understand the effect of disease on the population, it would serve as a warning for a presently unaffected population and device the vaccination strategies better for disease control. In this pursuit, we trained and optimized ML models with state of the art techniques and tested them on unseen artificial and real world networks. The results for SVR with RBF kernel and ANN show that $R_0$ was accurately predicted for these networks based on their known network properties. Hence, we have an estimate of disease controllability beforehand, without the need to simulate the SIR model on the test network.

## 3 METHODOLOGY

In this section we describe the procedure for generation of simulated data set. We also briefly describe the k-fold validation which a standard procedure used in ML model training and testing in the SI.

### 3.1 Generation of simulated data set

For simulation of the SIR model on networks, the parameters related to disease were fixed at: $k = 0.1$, $p_{ir} = 60\%$, $p_{id} = 30\%$, $p_{rs} = 10\%$. The starting population of individuals in each of the states was fixed at $S_o = 99.5\%$, $I_o = 0.5\%$ and $R_o = 0\%$. Each network had 1000 nodes and the network structure remains fixed throughout the simulation. The simulations were performed for 100-time steps and parameters $a$, $b$, $c$ and $e$ were determined using equations 2-5 in SI respectively. $R_0$ was calculated (using these parameters) and averaged over last 20 time steps, where the system reaches a stable regime (Fig.1 of SI). The networks were obtained using the python library NetworkX [Hagberg et al.(2008)Hagberg, Swart, and S Chult], which returns a network as output, for the supplied input parameter(s) governing connectivity patterns. For obtaining n (say) number of networks of a model network type, n values of these parameters were chosen from a range. This range was carefully chosen, such that the network properties fall in more or less the same range for all the models. Around 500 networks of each model kind (exact number in the description below) each of size($N$) 1000 were obtained. We use five model networks: $Erdos - Rényi$ (ER) random networks [ERDdS and R&WI(1959)], *Watts-Strogatz* small world (SW) networks [Watts and Strogatz(1998)], Scale Free (SF) networks [Bollobás et al.(2003)Bollobás, Borgs, Chayes, and Riordan], *Barabasi-Albert* [Barabási and Albert(1999)] (BA) and Stochastic block model (SBM) networks [Tcoyze(2016)] in our work. The parameters and their range of variation for each model network are described in the **Generation of model network examples** section in SI. The six network structural properties that were used as input features for training were Average Degree (avgdeg), Average Shortest Path Length (spl), Clustering Coefficient (cc), Network Density (den), Network Diameter (dia) and Maximum Degree (maxdeg). The definitions of these network metrics are presented in section **Complex Networks: Types and properties** of SI.

### 3.2 Training data-set and $k$-fold validation

For training the model, the *k-fold cross validation* routine of the sklearn library [Pedregosa et al.(2011)Pedregosa, Varoquaux, Gramfort, Michel, Thirion, Grisel, Blondel, Prettenhofer, Weiss, Dubourg, Vanderplas, Passos, Cournapeau, Brucher, Perrot, and Duchesnay] was used. The *k*-fold cross-validation (CV) procedure avoids over-fitting by holding back a part of data from use in training the model; such that the model performance is evaluated and reported based on testing on the unseen data. The full data set is split into k folds, out of which $k - 1$ folds are used for training the model and the remaining 1 fold is used for testing the model performance. The model performance

score is recorded every time, and then the model is discarded. This procedure is carried out in a loop ($k$ times) with a different fold held out for testing and $k - 1$ folds used for training every time. In this way, each fold is used once for testing and $k - 1$ times for training. Hence the model accuracy is averaged over all iterations of the procedure. In the present work all the ML model performances are evaluated using 10-fold cross validation technique. Also, please note that the data matrix was always permuted over rows before splitting it into testing and training parts, so that same model networks are not stacked together.

## 4 MODEL PERFORMANCE

The linear regression model resulted in a good fit when the networks from the same model were used for training and testing. The MSE and $R^2$ scores for the $ER$, $SF$, $SW$, $BA$ and $SBM$ networks as : (0.01, 0.99), (0.14, 0.87), (0.02, 0.99), (0.0, 0.99), (0.1, 0.99) respectively, indicating a good fit (correspondong plots shown in SI2). However, the linear regression lost the accuracy significantly when networks from all the models were used for training and testing (see Figure SI2) with MSE and $R^2$ as (4.99, 0.69). This shows that linear-regression is not a reliable model for predicting $R_0$ irrespective of the network type. Also, the failure of linear regression confirms the absence of linear relation between the input and target variable and calls for testing of non-linear regression techniques.

### 4.1 Support Vector Regression

Owing to the failure of linear regression in accurately predicting $R_0$, we explored the performance of SVR with RBF kernel on the data set. The performances for the other two kernels (linear, polynomial) have also been reported (see Table 1 for results). The parameters for polynomial and RBF kernel functions are $\gamma = 0.1$, $degree(d) = 2$ and $\gamma = \frac{1}{no.\ of\ features}$ (refer Eqns:13-14 in SI) respectively. SVR with linear,

TABLE 1
Table of Model performance results

| Model | Description | (MSE, $R^2$) |
|---|---|---|
| LR | | (4.99, 0.69) |
| SVR | Linear Kernel | (3.67, 0.73) |
| | Polynomial kernel | (11.30, 0.16) |
| | RBF kernel | (0.01, 1.00) |
| ANN | | (0.093, 0.99) |

polynomial and RBF kernels show mean squared error and $R^2$ as (3.67, 0.73), (11.30, 0.16) and (0.01, 1.00). Increasing the degree of the polynomial kernel to 3 improves the accuracy scores (3.02, 0.81); increasing the degree beyond 3 resulted in an arbitrarily high error

and requires much higher model training time than for degree 2. For the RBF kernel, the previously mentioned parameters were optimal concerning the accuracy and the required training time. Comparing the accuracy scores, we found that RBF kernel outperforms the other two kernels with a substantial margin and hence RBF can be used to predict $R_0$ with good precision. Please refer to Figure. 1 top panel and bottom panel (left figure) for the match of predicted output with the expected values, for all three kernels in SVR. The plots show predictions on only the first hundred data samples for better visualization.
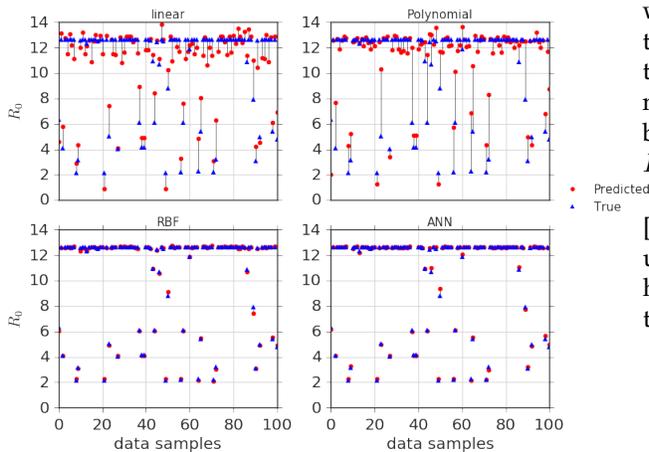


Fig. 1. The figure shows the difference in predicted and true $R_0$ using vertical lines at data points for linear, polynomial, RBF kernel in SVR and ANN respectively. For better visualization of results only $100$ randomly selected data points(true and corresponding predicted) are shown for all the models.

## 4.2 Neural Network Model

We also use a NN architecture (see Figure 2, left panel) that is optimized iteratively to gain maximum accuracy for regression. The NN model used here consists of three layers. The number of neurons in the input layer are conventionally fixed to be equal to the number of features of the data matrix. The output layer has one neuron, as the model performs regression to predict a number as an output ($R_0$). The hidden layer has 23 neurons that gather information from each neuron in the input layer and transmit it to the output layer. The number of neurons in the hidden layer were determined according to an empirical rule of thumb [hobs (https://stats.stackexchange.com/users/15974/hobs)()] that puts an upper limit on the total number of neurons without incurring the problem of over-fitting. The rule is,

$$N_h = \frac{N_s}{(\alpha \, (N_i + N_o))} \quad (1)$$

where $N_i$ is the number of input neurons, $N_o$ is the number of output neurons, $N_s$ is the number of samples in the training data set, and $\alpha$ is an arbitrary scaling factor between 2 and 10. For $\alpha = 2$, we get the maximum number of neurons according to the above formula, and any number of neurons greater than this value will result in over-fitting. For our case, we chose $\alpha = 10$, to avoid over-fitting and reduce the number of free parameters or weights in the model. Putting the known values of $N_i$, $N_o$ and $N_s$ as 6, 1 and 2552, we obtained $N_h = 36$. The optimal value of $N_h$ was then evaluated numerically by varying $N_h$ over a range of numbers within 36. The accuracy metrics were evaluated for a different number of neurons in the hidden layer, and this exercise was repeated for ten trials on randomly permuted data set. The optimum number of neurons in hidden layer were found out to be 23, as can be seen from the variation of MSE and $R^2$ coefficient in Figure. 2, right panel.

We used Keras (deep learning library for Python) [Chollet et al.(2015)] for model construction and simulation. Other specifics of the model in terms of its hyper-parameters and parameters are as described in the Table 1 in SI. The weights of edges in the NN
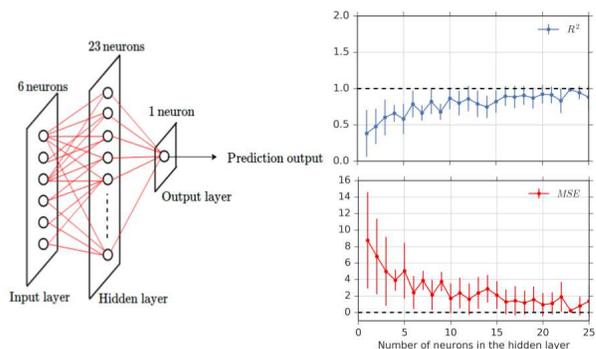


Fig. 2. The left panel shows a schematic of the NN model used in the present work; with 3 layers and number of neurons in each layer specified. The figures in the right panel show $R^2$ coefficient (top) and MSE (bottom) for $R_0$ prediction, averaged over 10 realizations for different number of neurons in the hidden layer; it can be seen that for $23$ number of neurons MSE touches the zero mark and $R^2$ touches the value one. This implies that using $23$ number of neurons in the hidden layer gives the maximum accuracy.

network were chosen from a normal distribution using the kernel initializer function. The activation functions for neurons were governed by the rectified linear unit ("relu") i.e., the neuron activation was linearly related to the input. Adaptive Moment Estimation ("Adam") was used as optimizer, which is based on an adaptive learning scheme for updating of the network weights. This optimizer function updates the learning rates iteratively based on the moments of the gradient of the objective function. The objective function or the

loss function was accuracy measured in terms of MSE. With epoch size ("Epochs") set at 50, the batch size of 5 and other parameters set as specified above in the NN model, the mean accuracy measured in terms of (MSE, $R^2$) for 10-fold cross validation was $(0.093, 0.99)$.

We have shown the predicted and true $R_0$ for all the examples using SVR (with RBF kernel) and ANN model in Figure 4. We also trained all the ML model using only four (*avgdeg*, *maxdeg*, *dia*, *spl*) out of six features selected based on their contribution indices. These four features had highest values of the contribution indices (refer to subsection **Ranking of the features** in SI). We have shown the relative contribution indices for all the six features in the Figure 3. We observe that model performances are still fairly accurate in predicting the correct $R_0$. In the Table 2, the model performance metrics for SVR with RBF and ANN for training with six and four features respectively have been shown. We can infer from the results that there is a trade-off between accuracy of model predictions and number of features used for training the model. The accuracy of the predicted value is better with the all the six features in the data set than when only four feature vectors were considered. The precision of the prediction with top-four features is slightly reduced but it is in a bearable range (refer to the MSE and variance scores in table2). Therefore, one can remove *cc*, *den* from the feature set for the training without compromising much with the prediction accuracy.

| Number of Features | ML technique | Accuracy Measures | |
|---|---|---|---|
| | | MSE | $R^2$ |
| Four | SVR(RBF) | 0.11 | 0.99 |
| | ANN | 2.99 | 0.82 |
| Six | SVR(RBF) | 0.01 | 1.00 |
| | ANN | 0.013 | 0.998 |

TABLE 2
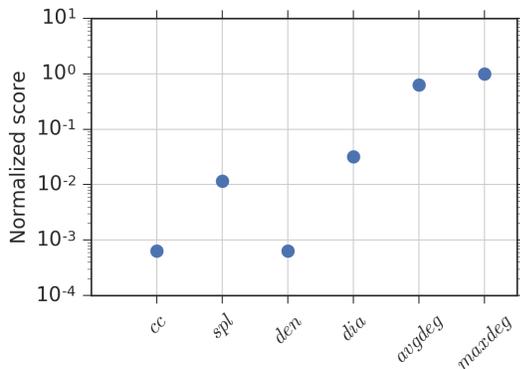Table of ML model performance results on simulated networks.



Fig. 3. This figure shows relative importance of the network features based on their contribution index. The contribution indices are normalized between 0 and 1
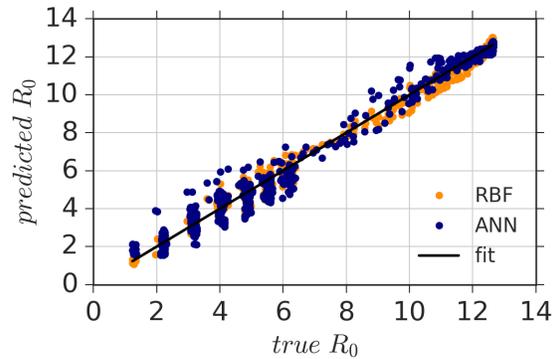


Fig. 4. The figure shows the predicted $R_0$ vs true $R_0$ for ANN model and SVR model with RBF kernel for all the data points. The black line (fit) corresponds to ideal case where true $R_0$ is equal predicted $R_0$.

### 4.3 Performance on Real-world Networks

We tested the ML models for $R_0$ prediction for four real-world network datasets: *infect-dublin* [Isella et al.(2011)Isella, Stehlé, Barrat, Cattuto, Pinton, and Van den Broeck], *infect-hyper* [Isella et al.(2011)Isella, Stehlé, Barrat, Cattuto, Pinton, and Van den Broeck], *crime-moreno* and *email-univ* ( [Rossi and Ahmed(2015)], [Guimera et al.(2003)Guimera, Danon, Diaz-Guilera, Giralt, and Arenas]). *infect-dublin* and *infect-hyper* are categorized as proximity networks based on face-to-face interactions between people, with the number of nodes ($N$) and number of edges (E) being $(410, 2765)$ and $(113, 2196)$ respectively. *crime-moreno* is categorized as interaction network with ($N$, $E$) = $(829, 1474)$. *email-univ* is email network with ($N$, $E$) = $(1133, 5451)$. Please note that we chose networks with single giant component only.

The accuracy metrics corresponding to the ML models on real-world networks is tabulated in Table 3. It is observed that the ML model performs very accurately for these networks as well, as for artificial networks. Especially, for *infect-dublin* and *crime-moreno* networks, both SVR and ANN models predict $R_0$ that almost matches the true $R_0$ value. Furthermore, these real-world test networks serve as unseen test examples for the ML models, and accurate prediction of $R_0$ authenticates the ML models even more.

## 5 CONCLUSION AND FUTURE PROSPECTS

The present work explored the applicability of ML regression techniques to predict basic reproduction number, $R_0$, a factor indicative of the effect of disease or its controllability, on complex networks. $R_0$, in general, depends on many factors: the duration of disease persistence in the population, the vulnerability of an individual to an infection, the number of infected neighbours to a susceptible individual, etc.

TABLE 3
Table of Model performance results on real world networks

| Dataset | True $R_o$ | Pred. $R_o$ (SVR) | Pred. $R_o$ (ANN) |
|---|---|---|---|
| infect-dublin | 5.63 | 5.97 | 5.13 |
| infect-hyper | 10.02 | 8.09 | 10.27 |
| crime-moreno | 1.95 | 1.75 | 2.01 |
| email-univ | 4.22 | 4.19 | 3.67 |

On the other hand, if we have a population where all these above parameters are fixed to a reasonable value, how the social strata(complex network in our case) on which the disease spreads affect the disease spreading is still a question. To explore this, we examined whether $R_0$ can be predicted based on global properties of the network in hand, irrespective of the model network type it belongs to?

A large number of networks were generated, and dynamics of the disease spreading were simulated on these networks, and the corresponding $R_0$ was recorded along with the network properties. Using the recorded data, three ML regression models were trained to predict $R_0$ values on the test data. These models were tuned based on their parameters to obtain good prediction. The results using RBF kernel in SVR and ANN models showed high accuracy of $R_0$ prediction, suggesting that there exists a significant correlation between the network properties and disease controllability. The generalizability of the trained models is convincing because the testing was always performed on unseen data using the k-fold validation technique. One of the improvements to the present work could be to train the models using a larger number networks of different sizes, such that a higher range of network properties such as shortest path length, clustering coefficient, etc. is spanned. These models will then yield correct prediction for any given test network. However, as one can see, this is just a scalability issue. Moreover, good predictions of these models on real-world networks is an exciting result. Our work reports two significant findings (a) The disease controllability on the network can be predicted using global network properties. (b) The standard ML techniques can be applied to processes on complex network. In our case it is predicting disease controllability on a network. The tunability of ML models offers immense power to forecast or predict processes on complex network systems.

The computational cost for some of the features is high (especially for the clustering coefficient (*cc*) and the shortest path length (*spl*)). Hence, the time complexity for obtaining the features for the training data set will be high. This is one of the constraints of our approach for the ML model training. But, for predicting the value of $R_0$ for any arbitrary test network based on known network features, the prediction time is almost negligible. This implies that we can predict the value of $R_0$ at the very first stage of getting the test network (without waiting until the epidemic outbreak has completed or reached a stable state). Of course, underlying assumption is that we should know the value of these features beforehand at the time of testing.

The prospects of the work may include using deep learning approaches for unsupervised learning of features. As we know that the numerical calculation of network properties for the training as well as testing the model is a time-consuming step, it would be great if the network itself could be made to train the model. Another prospect is to explore if the network adjacency matrix can be used to train a deep learning CNN architecture. Also, if network embedding algorithms can be used to learn the features that are instrumental in the disease spreading, it would be a significant leap forward from this work.

## REFERENCES

[Anderson and May(1992)] R. M. Anderson and R. M. May, *Infectious diseases of humans: dynamics and control.* Oxford university press, 1992.

[Barabási and Albert(1999)] A.-L. Barabási and R. Albert, "Emergence of scaling in random networks," *science*, vol. 286, no. 5439, pp. 509–512, 1999.

[Hethcote(2000)] H. W. Hethcote, "The mathematics of infectious diseases," *SIAM review*, vol. 42, no. 4, pp. 599–653, 2000.

[Lotka(1956)] A. J. Lotka, "Elements of mathematical biology," 1956.

[Stewart et al.(2005)Stewart, Logsdon, and Kelley] A. D. Stewart, J. M. Logsdon, and S. E. Kelley, "An empirical study of the evolution of virulence under both horizontal and vertical transmission," *Evolution*, vol. 59, no. 4, pp. 730–739, 2005.

[Shirley and Rushton(2005)] M. D. Shirley and S. P. Rushton, "The impacts of network topology on disease spread," *Ecological Complexity*, vol. 2, no. 3, pp. 287–299, 2005.

[Schimit and Pereira(2018)] P. H. Schimit and F. H. Pereira, "Disease spreading in complex networks: A numerical study with principal component analysis," *Expert Systems with Applications*, vol. 97, pp. 41–50, 2018.

[MacQueen et al.(1967)] J. MacQueen *et al.*, "Some methods for classification and analysis of multivariate observations," in *Proceedings of the fifth Berkeley symposium on mathematical statistics and probability*, vol. 1, no. 14. Oakland, CA, USA, 1967, pp. 281–297.

[Chaplot et al.(2006)Chaplot, Patnaik, and Jagannathan] S. Chaplot, L. Patnaik, and N. Jagannathan, "Classification of magnetic resonance brain images using wavelets as input to support vector machine and neural network," *Biomedical signal processing and control*, vol. 1, no. 1, pp. 86–92, 2006.

[Xin et al.(2018)Xin, Zhang, and Shao] R. Xin, J. Zhang, and Y. Shao, "Complex network classification with convolutional neural network," *arXiv preprint arXiv:1802.00539*, 2018.

[Pedregosa et al.(2011)Pedregosa, Varoquaux, Gramfort, Michel, Thirion, Grisel, Blor F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau,

M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

[Hagberg et al.(2008)Hagberg, Swart, and S Chult] A. Hagberg, P. Swart, and D. S Chult, "Exploring network structure, dynamics, and function using networkx," Los Alamos National Lab.(LANL), Los Alamos, NM (United States), Tech. Rep., 2008.

[ERDdS and R&WI(1959)] P. ERDdS and A. R&WI, "On random graphs i," *Publ. Math. Debrecen*, vol. 6, pp. 290–297, 1959.

[Watts and Strogatz(1998)] D. J. Watts and S. H. Strogatz, "Collective dynamics of small-worldnetworks," *nature*, vol. 393, no. 6684, p. 440, 1998.

[Bollobás et al.(2003)Bollobás, Borgs, Chayes, and Riordan] B. Bollobás, C. Borgs, J. Chayes, and O. Riordan, "Proceedings of the fourteenth annual acm-siam symposium on discrete algorithms," 2003.

[Tcoyze(2016)] Tcoyze, "Project title," Mar. 2016. [Online]. Available: https://github.com/tcoyze/stochastic-blockmodel

[hobs (https://stats.stackexchange.com/users/15974/hobs)()] hobs (https://stats.stackexchange.com/users/15974/hobs), "How to choose the number of hidden layers and nodes in a feedforward neural network?" Cross Validated, uRL:https://stats.stackexchange.com/q/136542 (version: 2017-04-13). [Online]. Available: https://stats.stackexchange.com/q/136542

[Chollet et al.(2015)] F. Chollet *et al.*, "Keras," https://keras.io, 2015.

[Isella et al.(2011)Isella, Stehlé, Barrat, Cattuto, Pinton, and Van den Broeck] L. Isella, J. Stehlé, A. Barrat, C. Cattuto, J.-F. Pinton, and W. Van den Broeck, "What's in a crowd? analysis of face-to-face behavioral networks," *Journal of theoretical biology*, vol. 271, no. 1, pp. 166–180, 2011.

[Rossi and Ahmed(2015)] R. A. Rossi and N. K. Ahmed, "The network data repository with interactive graph analytics and visualization," in *AAAI*, 2015. [Online]. Available: http://networkrepository.com

[Guimera et al.(2003)Guimera, Danon, Diaz-Guilera, Giralt, and Arenas] R. Guimera, L. Danon, A. Diaz-Guilera, F. Giralt, and A. Arenas, "Self-similar community structure in a network of human interactions," *Physical review E*, vol. 68, no. 6, p. 065103, 2003.