Capturing the Effects of Transportation on the Spread of COVID-19 with a Multi-Networked SEIR Model

Damir Vrabac, Mingfeng Shang, Brooks Butler, Joseph Pham, Raphael Stern, Philip E. Paré

Abstract—In this paper we present a deterministic discrete-time networked SEIR model that includes a number of transportation networks, and present assumptions under which it is well defined. We analyze the limiting behavior of the model and present necessary and sufficient conditions for estimating the spreading parameters from data. We illustrate these results via simulation and with real COVID-19 data from the Northeast United States, integrating transportation data into the results.

Index Terms—Control applications, transportation networks, SEIR model, COVID-19

I. INTRODUCTION

N December 2019, a novel coronavirus (SARS-CoV-2), that causes the disease COVID-19, was detected in Wuhan, China. This virus quickly spread throughout China, and before long, the virus had reached the status of a global pandemic. In order to minimize the impact of COVID-19, it is critical to be able to quickly track the spread of the virus and understand the mechanisms that are enabling its propagation. While the mode of transmission of the virus is not exactly known, human-to-human interaction appears to be a main factor [1]. A key component for transmission is the underlying transportation network, which acts as a propagator of the virus within and between communities.

In this work we extend the deterministic SEIR [2] model for viral spread to consider spread over the network in the context of human interaction and transportation. We model the proportion of people in each county who have not been infected (S), those who have been infected but have not been confirmed via a test (E), test-confirmed infected cases (I), and those who have either recovered or died from the virus (R) and show that we are able to model the evolution of such a virus, as well as recover the proper model parameter values from time series data of infections and recoveries and apply this model to the recent COVID-19 outbreak.

This manuscript was first submitted for review on September 14, 2020. This material is based upon work supported by the National Science Foundation under Grant No. CNS-2028946 (R.S.) and Grant No. CNS-2028738 (P.E.P.).

Damir Vrabac is with the Department of Computer Science at Stanford University. Mingfeng Shang, Joseph Pham, and Raphael Stern are with the Department of Civil, Environmental, and Geo-Engineering at the University of Minnesota. Brooks Butler and Philip E. Paré are with the School of Electrical and Computer Engineering at Purdue University. E-mails: dvrabac@stanford.edu, shang140@umn.edu, brooksbutler@purdue.edu, pham0231@umn.edu, rstern@umn.edu, philpare@purdue.edu.

The SEIR model has become popular for modeling epidemic spread (e.g., [3]) and has been described in [4]. A similar model to the SEIR model, the SEIV model has been studied in previous work where the vigilant state V corresponds to a state that is not infected nor immediately susceptible, i.e. similar to the recovered state R [5]. The model has also been extended to account for quarantine [6] and asymptomatic transmission [7]. When considering how transportation can propagate a viral outbreak, the SIS model has been extended to include transportation flows between nodes [8]. We go beyond prior work by integrating transportation networks into a networked SEIR model, analyzing the model, and applying it to the COVID-19 pandemic.

The multi-networked SEIR model is introduced in Section II and its limiting behavior is discussed in Section III. Results on model parameter estimation are given in Section IV, and the model is applied via simulations and real COVID-19 data to the Northeast US. We conclude in Section VI.

A. Notation

Given a vector x, the transpose is indicated by x^{\top} , \bar{x} is the average of its entries, and $\operatorname{diag}(\cdot)$ is a diagonal matrix with the argument on the diagonal. We use $\mathbf{0}$ and $\mathbf{1}$ to denote a vector or matrix of zeros and ones, respectively, of the appropriate dimensions. We define a directed graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, w)$, where \mathcal{V} is the set of nodes, $\mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$ is the set of edges, and $w: \mathcal{E} \to \mathbb{R}^+$ is a function mapping directed edges to their weightings, with \mathbb{R}^+ being the set of positive real values. Given \mathcal{G} , we denote an edge from node $i \in \mathcal{V}$ to node $j \in \mathcal{V}$ by (i,j). We say node $i \in \mathcal{V}$ is a neighbor of node $j \in \mathcal{V}$ if and only if $(i,j) \in \mathcal{E}$, and denote the neighbors of node j as \mathcal{N}_j . We denote the weighted adjacency matrix associated with \mathcal{G} as A with the nonzero entry a_{ji} indicating the strength of edge (i,j) as given by w. We use [n] to denote the set $\{1,2,...,n\}$.

II. MULTI-NETWORKED SEIR MODEL

Here we introduce the discrete-time multi-networked SEIR model. We assume that the virus spreads over a set of graphs $\mathcal{G}^l = (\mathcal{V}, \mathcal{E}^l, w^l)$, for $l \in \mathcal{L}$, where we interpret each node in \mathcal{V} as a subpopulation, \mathcal{L} is the set of transportation networks, and \mathcal{E}^l, w^l capture the weighted transportation links. Node i's susceptibility s_i^k , exposed e_i^k , infection p_i^k , and recovery r_i^k

proportions evolve as

$$s_i^{k+1} = s_i^k - h s_i^k \iota_i^k, \tag{1a}$$

$$e_i^{k+1} = e_i^k + h s_i^k \iota_i^k - h \sigma_i e_i^k,$$
 (1b)

$$p_i^{k+1} = p_i^k + h(\sigma_i e_i^k - \gamma_i p_i^k), \tag{1c}$$

$$r_i^{k+1} = r_i^k + h\left(\gamma_i p_i^k\right),\tag{1d}$$

where k is the time step, h is the sampling parameter, σ_i captures the rate at which the exposed become confirmed infected cases, γ_i is the recovery rate, and

$$\iota_i^k = \sum_{l \in \mathcal{L}} \left(\beta_i^{e,l} \sum_{j \in \mathcal{N}_i^l} a_{ij}^l e_j^k + \beta_i^{p,l} \sum_{j \in \mathcal{N}_i^l} a_{ij}^l p_j^k \right), \quad (2)$$

where a_{ij}^l represents the edge weights, and $\beta_i^{e,l}$ and $\beta_i^{p,l}$ are the corresponding infection rates for the lth transportation network. Note for the special case where $|\mathcal{L}|=1$, the model in (1)-(2) becomes the traditional networked SEIR model.

For the discrete-time SEIR model to be well-defined we need the following assumption.

Assumption 1: For all $i \in [n]$, we have $0 < h\gamma_i < 1$, $\begin{array}{l} 0 < h\sigma_i < 1, \ 0 \leq h\check{\iota}_i^k < 1, \ \text{where} \ \ \check{\iota}_i^k = \sum_{l \in \mathcal{L}} (\beta_i^{e,l} + \beta_i^{p,l}) \sum_{j \in \mathcal{N}_i} a_{ij}^l, \ \text{and} \ \beta_i^{e,l}, \beta_i^{p,l}, a_{ij}^l \geq 0, \ \text{for all} \ j \in [n]. \end{array}$ Assumption 1 requires the sampling parameter to be small enough in relation to the healing parameters and the denseness of the graph scaled by the infection parameters, and guarantees that the model is well defined.

Lemma 1: Consider the model in (1)-(2) under Assumption 1. Suppose $s_i^0, e_i^0, p_i^0, r_i^0 \in [0, 1], s_i^0 + e_i^0 + p_i^0 + r_i^0 = 1$ for all $i \in [n]$. Then, for all $k \ge 0$ and $i \in [n]$, $s_i^k, e_i^k, p_i^k, r_i^k \in [0, 1]$ and $s_i^k + e_i^k + p_i^k + r_i^k = 1$.

Proof: We prove this result by induction. By assumption, it holds for the base-case k = 0. We follow the proof by showing the induction-step, that is, assume $s_i^k, e_i^k, p_i^k, r_i^k \in$ [0,1] and $s_i^k + e_i^k + p_i^k + r_i^k = 1$, for all $i \in [n]$, and we now show that this holds also for time-step k+1. By Assumption 1 and (1a), $s^{k+1} \ge s^k_i - h s^k_i \check{\iota}^k_i = s^k_i \left[1 - h \check{\iota}^k_i\right] \ge 0$. We also have $s^{k+1} \le s^k \le 1$ since $h\left[-s^k_i \iota^k_i\right] \le 0$. By Assumption 1 and (1b), $e_i^{k+1} \geq (1 - h\sigma_i)e_i^k \geq 0$. Moreover, by the assumption $e_j^k, p_j^k \leq 1$ for all $j \in [n]$, Assumption 1, and (1b), $e_i^{k+1} \leq 1$ $\begin{array}{l} c_i^k + s_i^k h \check{t}_i^k \leq e_i^k + s_i^k \leq 1. \text{ By Assumption 1 and (1c), } p_i^{k+1} \geq \\ (1-h\gamma_i) \, p_i^k \geq 0 \text{ and } p_i^{k+1} \leq p_i^k + h\sigma_i e_i^k \leq p_i^k + e_i^k \leq 1. \text{ By Assumption 1 and (1d), } r_i^{k+1} \geq r_i^k \geq 0, \text{ and } r_i^{k+1} \leq r_i^k + p_i^k. \end{array}$

Thus, by the principle of mathematical induction we have that, if $s_i^0, e_i^0, p_i^0, r_i^0 \in [0,1]$ and $s_i^0 + e_i^0 + p_i^0 + r_i^0 = 1$ for all $i \in [n]$ then $s_i^k, e_i^k, p_i^k, r_i^k \in [0, 1]$ and $s_i^k + e_i^k + p_i^k + r_i^k = 1$ for all $k \in \mathbb{N}$.

III. ANALYSIS OF MODEL

In this section we present a result on the stability of the healthy states of the networked SEIR model, that is, where $\lim_{k\to\infty} e_i^k = 0$ and $\lim_{k\to\infty} p_i^k = 0$ for all $i\in[n]$.

Let $\lambda_{max}^{M_k}$ be the dominant eigenvalue of M_k , where M_k is defined as

$$M_k = \begin{bmatrix} (I + hS^kT^e - h\sigma) & hS^kT^p \\ h\sigma & (I - h\gamma) \end{bmatrix}, \tag{3}$$

where $S^k = \operatorname{diag}(s_i^k)$, $T^e = \sum_{l \in \mathcal{L}} B_l^e A_l$, $T^p = \sum_{l \in \mathcal{L}} B_l^p A_l$, $B_l^e = \operatorname{diag}(\beta_i^{e,l})$, $B_l^p = \operatorname{diag}(\beta_i^{p,l})$, $\gamma = \operatorname{diag}(\gamma_i)$ and $\sigma = \sum_{l \in \mathcal{L}} B_l^p A_l$ $\operatorname{diag}(\sigma_i)$. Note that M_k captures the dynamics of the vector of the exposed and infection states, e^k and p^k , in (1)-(2).

Theorem 1: Consider the model in (1)-(2) under Assumption 1. Suppose $s_i^0, e_i^0, p_i^0, r_i^0 \in [0, 1], s_i^0 + e_i^0 + p_i^0 + r_i^0 = 1$ for all $i \in [n]$, T^p is irreducible, $s_i^0 > 0$ for all $i \in [n]$, and $p_i^0 > 0$ for some i. Then, for all $k \ge 0$ and $i \in [n]$,

- 1) $s_i^{k+1} \le s_i^k$,
- 2) $\lim_{k\to\infty} e_i^k = 0$ and $\lim_{k\to\infty} p_i^k = 0$,
- 3) $\lambda_{max}^{M_k}$ is monotonically decreasing as a function of k,
- 4) there exist a \bar{k} such that $\lambda_{max}^{M_k} < 1$ for all $k \geq \bar{k}$, 5) there exists \bar{k} , such that p_i^k converges linearly to 0 for all $k > \bar{k}$ and $i \in [n]$.

Proof: We present the proof for each part of the theorem, starting with 1).

- 1) By Lemma 1 and Assumption 1, we have that $-hs_i^k \iota_i^k \leq$ 0 for all $i \in [n]$ and $k \ge 0$. Therefore, from (1a), we have $s_i^{k+1} \le s_i^k.$
- 2) Since the rate of change of s^k , $-hS^k \left[T^e e^k + T^p p^k \right]$, is non-positive for all $k \geq 0$ and s^k is lower bounded by zero, by Lemma 1, we conclude that $\lim_{k\to\infty} s^k$ exists. Therefore,

$$\lim_{k \to \infty} -hS^k \left[T^e e^k + T^p p^k \right] = \mathbf{0}. \tag{4}$$

Therefore, $\lim_{k\to\infty} e^{k+1} - e^k = \lim_{k\to\infty} -h\sigma e^k$. Thus, by Assumption 1, $h\sigma_i > 0$ for all $i \in [n]$, $\lim_{k \to \infty} e_i^k = 0$ for all $i \in [n]$.

Similarly, we show that $\lim_{k\to\infty} p_i^k = 0$ for all $i \in [n]$. We have that $\lim_{k\to\infty} p^{k+1} - p^k = \lim_{k\to\infty} h\left(\sigma e^k - \gamma p^k\right) =$ $\lim_{k\to\infty} -h\gamma p^k$, where we used that $\lim_{k\to\infty} e^k = 0$. By assumption $h\gamma_i > 0$ for all $i \in [n]$, thus $\lim_{k \to \infty} p_i^k = 0$ for all $i \in [n]$.

- 3) By assumption $s_i^0 > 0$ for all $i \in [n]$, and from the proof of Lemma 1 we can see that $s_i^k > 0$ for all $i \in [n]$, $k \geq 0$. Therefore, since we have that T^p is irreducible, from (3) and Assumption 1, the matrix M_k is irreducible and non-negative, for all finite k. Thus by the Perron-Frobenius Theorem for irreducible non-negative matrices we have that $\lambda_{max}^{M_k} = \rho(M_k)$. Since $\rho(M_k)$ increases when any entry increases [9, Theorem 2.7] and by 1) of this theorem, we have that $\rho(M_k) \ge \rho(M_{k+1})$, that is $\lambda_{max}^{M_k} \ge \lambda_{max}^{M_{k+1}}$.
- 4) There are two possible equilibria: i) $\lim_{k\to\infty} s^k = \mathbf{0}$, and ii) $\lim_{k\to\infty} s^k = s^* \neq \mathbf{0}$. We explore the two cases separately.
 - i) If $\lim_{k\to\infty} s^k = \mathbf{0}$,

$$\lim_{k\to\infty} M_k = \begin{bmatrix} I-h\sigma & 0\\ h\sigma & I-h\gamma \end{bmatrix}.$$

Therefore, by Assumption 1, there exists a \bar{k} such that $\lambda_{max}^{M_k}$

ii) If $\lim_{k\to\infty} s^k = s^* \neq \mathbf{0}$, then, by 2), for any (s^0, e^0, p^0, r^0) the system converges to some equilibrium of the form $(s^*, 0, 0, 1 - s^*)$. Define

$$\epsilon_s^k := s^k - s^* \text{ and } \epsilon_p^k := \begin{bmatrix} e^k \\ p^k \end{bmatrix} - \mathbf{0}_{2n}.$$
 (5)

Linearizing the dynamics of ϵ_s^k and ϵ_p^k around $(s^*, \mathbf{0}_{2n})$ gives

$$\epsilon_s^{k+1} = \epsilon_s^k - hS^k \begin{bmatrix} T^e & T^p \end{bmatrix} \epsilon_p^k, \tag{6a}$$

$$\epsilon_n^{k+1} = M_k \epsilon_n^k. \tag{6b}$$

Let
$$\lambda_{max}^{M^*}$$
 be the maximum eigenvalue of
$$M^* = \begin{bmatrix} (I + h \operatorname{diag}(s^*) T^e - h\sigma) & h \operatorname{diag}(s^*) T^p \\ h\sigma & (I - h\gamma) \end{bmatrix} \quad (7)$$

$$w^{*\top}M^{*} = \lambda_{max}^{M^{*}}w^{*\top}.$$
 (8)

If $\lambda_{max}^{M^*} > 1$, then the system in (6) is unstable. Therefore, by Lyapunov's Indirect Method, $\lim_{k \to \infty} \left(\epsilon_s^k, \epsilon_p^k \right) \neq (s^*, \mathbf{0}_{2n})$, which is a contradiction.

Now consider the case where $\lambda_{max}^{M^*} = 1$. Define $\tilde{M}_k = \begin{bmatrix} h \mathrm{diag}(\epsilon_s^k) T^e & h \mathrm{diag}(\epsilon_s^k) T^p \\ 0 & 0 \end{bmatrix}$.

Then we can write $M_k=M^*+\tilde{M}_k$, observe that all entries in \tilde{M} are non-negative. Using (5) and left multiplying the equation of ϵ_p^{k+1} in (6b) by $w^{*\top}$ we get

$$w^{*\top} \epsilon_p^{k+1} = w^{*\top} M_k \epsilon_p^k$$

$$= \lambda_{max}^{M^*} w^{*\top} \epsilon_p^k + w^{*\top} \tilde{M}_k \epsilon_p^k$$

$$= w^{*\top} \epsilon_p^k + w^{*\top} \tilde{M}_k \epsilon_p^k.$$

Thus,

$$w^{*\top} \left(\epsilon_p^{k+1} - \epsilon_p^k \right) = w^{*\top} \tilde{M}_k \epsilon_p^k \ge 0, \tag{10}$$

where the last inequality holds since all elements are nonnegative. This contradicts that $\lim_{k\to\infty} z^k = \mathbf{0}_{2n}$, that is 2). Therefore, there exists a \bar{k} such that $\lambda_{max}^{M_k} < 1$ for all $k \geq \bar{k}$.

5) Since, by 4), there exists a \bar{k} such that $\lambda_{max}^{M_k} < 1$ for all

 $k \geq \bar{k}$, and we know that $\lambda_{max}^{M_k} = \rho(M_k) \geq 0$ by Assumption 1, we have

$$\lim_{k \to \infty} \frac{\|p^{k+1}\|}{\|p^k\|} = \frac{\|M_k p^k\|}{\|p^k\|} = \lambda_{max}^{M_k} < 1.$$
 (11)

Therefore, for $k > \bar{k}$, p^k converges linearly to $\mathbf{0}_n$. Note that the proof was inspired by a similar result for the SIR model [10, Theorem 1]. The results in Theorem 1 show that the virus will die out, providing insight into the convergence rate, under mild assumptions.

IV. ESTIMATING MODEL PARAMETERS

We now explore conditions for estimating the SEIR model parameters from data. Due to space limitations we consider $|\mathcal{L}|=1$ and refer to Remark 1 for $|\mathcal{L}|>1$. In order to estimate the parameters we define the following matrices:

$$\Phi = \begin{bmatrix} hS^{0}Ae^{0} & hS^{0}Ap^{0} & -he^{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ hS^{T-1}Ae^{T-1} & hS^{T-1}Ap^{T-1} & -he^{T-1} & \mathbf{0} \end{bmatrix}, (12)$$

$$\Sigma = \begin{bmatrix} 0 & 0 & he^0 & -hp^0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & he^{T-1} & -hp^{T-1} \end{bmatrix}, \tag{13}$$

and
$$\Gamma = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{d} \end{bmatrix}. \tag{14}$$

Using the above matrices we write (1) as

$$\begin{bmatrix} e^{1} - e^{0} \\ \vdots \\ e^{T} - e^{T-1} \\ p^{1} - p^{0} \\ \vdots \\ p^{T} - p^{T-1} \\ r^{1} - r^{0} \\ \vdots \\ r^{T} - r^{T-1} \end{bmatrix} = \underbrace{\begin{bmatrix} \Phi \\ \Sigma \\ \Gamma \end{bmatrix}}_{Q} \begin{bmatrix} \beta^{e} \\ \beta^{p} \\ \sigma \\ \gamma \end{bmatrix}. \tag{15}$$

We find the least squares estimates $\hat{\beta}^e$, $\hat{\beta}^p$, $\hat{\sigma}$, and $\hat{\gamma}$ using the pseudoinverse of Q.

Theorem 2: Consider the model in (1) with homogeneous virus spread, that is, β^e , β^p , σ , and γ are the same for all nnodes. Assume that s^k, e^k, p^k, r^k , for all $k \in [T] \cup \{0\}$, and h are known, with n > 1. Then, the parameters of the spreading process can be identified uniquely if and only if T > 0, and there exist $i_1, i_2, i_3, i_4 \in [n]$ and $k_1, k_2, k_3, k_4 \in [T-1] \cup \{0\}$ such that

$$p_{i_1}^{k_1} \neq 0, e_{i_2}^{k_2} \neq 0,$$

$$g_{i_2}^{k_3}(e^{k_3})g_{i_4}^{k_4}(p^{k_4}) \neq g_{i_4}^{k_4}(e^{k_4})g_{i_2}^{k_3}(p^{k_3}),$$
(16a)

$$g_{i_3}^{k_3}(e^{k_3})g_{i_4}^{k_4}(p^{k_4}) \neq g_{i_4}^{k_4}(e^{k_4})g_{i_3}^{k_3}(p^{k_3}),$$
 (16b)

where $g_i^k(x) = s_i^k \sum_{j \in \mathcal{N}_i} a_{ij} x_j$. *Proof*: Using (12)-(14), we can write Q as follows

$$Q = \underbrace{\begin{bmatrix} I & -I & \mathbf{0}_{nT \times nT} \\ \mathbf{0}_{nT \times nT} & I & -I \\ \mathbf{0}_{nT \times nT} & \mathbf{0}_{nT \times nT} & I \end{bmatrix}}_{\tilde{D}} \underbrace{\begin{bmatrix} \mathbf{a} & \mathbf{b} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{c} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{d} \end{bmatrix}}_{\tilde{Q}}. \quad (17)$$

Since n > 1, $\Phi = |\mathbf{a}| \mathbf{b}$ has at least two rows, and given that (16b) holds, $\tilde{\Phi}$ has column rank equal to two. Moreover, if (16a) holds c and d each have at least one element that is nonzero. Thus, \tilde{Q} has full column rank. Clearly D has full rank which implies that the rank of Q is equal to the rank of Q [11]. Therefore, there exists a unique solution to (15) using the pseudoinverse.

If one of the assumptions in (16a)-(16b) is not met, Q will have a nontrivial nullspace. Therefore, in that case, (15) does not have a unique solution.

For the heterogeneous case it is not necessary to know all entries of s^k, e^k, p^k, r^k . It is sufficient to know only $s^k_j, e^k_j, p^k_j, r^k_j$, for $j \in \mathcal{N}_{i_1} \cup \mathcal{N}_{i_2} \cup \mathcal{N}_{i_3} \cup \mathcal{N}_{i_4} \cup \{i_1, i_2, i_3, i_4\}$, where i_1, i_2, i_3, i_4 satisfy (16).

To estimate the spreading parameters for the discrete-time, heterogeneous SEIR model from Section II we define:

$$\Phi = \begin{bmatrix} hS^0Ae^0 & hS^0Ap^0 & -he^0 & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \\ hS^{T-1}Ae^{T-1} & \mathbf{b} \end{bmatrix}, \quad (12) \quad \Phi_i = \begin{bmatrix} hs_i^0\sum_{j\in\mathcal{N}_i}a_{ij}e_j^0 & hs_i^0\sum_{j\in\mathcal{N}_i}a_{ij}p_j^0 & -he_i^0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ hs_i^{T-1}\sum_{j\in\mathcal{N}_i}a_{ij}e_j^{T-1} & \vdots & \vdots & \vdots \\ hs_i^{T-1}\sum_{j\in\mathcal{N}_i}a_{ij}e_j^{T-1} & \underbrace{hs_i^{T-1}\sum_{j\in\mathcal{N}_i}a_{ij}p_j^{T-1}}_{\mathbf{b}_i} & -he_i^{T-1} & 0 \end{bmatrix}, \quad (12)$$

$$\Sigma_{i} = \begin{bmatrix} 0 & 0 & he_{i}^{0} & -hp_{i}^{0} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & he_{i}^{T-1} & -hp_{i}^{T-1} \end{bmatrix}, \tag{18}$$

 $\Gamma_i = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{d}_i \end{bmatrix}$. (20) Using the above matrices we write (1) as

$$\begin{bmatrix} e_{i}^{1} - e_{i}^{0} \\ \vdots \\ e_{i}^{T} - e_{i}^{T-1} \\ p_{i}^{1} - p_{i}^{0} \\ \vdots \\ p_{i}^{T} - p_{i}^{T-1} \\ r_{i}^{1} - r_{i}^{0} \\ \vdots \\ r_{i}^{T} - r_{i}^{T-1} \end{bmatrix} = \underbrace{\begin{bmatrix} \Phi_{i} \\ \Sigma_{i} \\ \Gamma_{i} \end{bmatrix}}_{Q_{i}} \begin{bmatrix} \beta_{i}^{e} \\ \beta_{i}^{p} \\ \sigma_{i} \\ \gamma_{i} \end{bmatrix}. \tag{21}$$

We find the least squares estimates $\hat{\beta}_i^e$, $\hat{\beta}_i^p$, $\hat{\sigma}_i$, and $\hat{\gamma}_i$ using the pseudoinverse of Q_i .

Corollary 1: Consider the model in (1). Assume that $s_i^k, e_j^k, p_j^k, r_i^k$, for all $j \in \mathcal{N}_i \cup \{i\}, k \in [T-1] \cup \{0\}, e_i^T, p_i^T, r_i^T$, and h are known. Then, the parameters of the spreading process for node i can be identified uniquely if and only if T>1, and there exist $k_1, k_2, k_3, k_4 \in [T-1] \cup \{0\}$ such that

$$\begin{aligned} p_i^{k_1} &\neq 0, e_i^{k_2} \neq 0, \\ g_i^{k_3}(e^{k_3}) g_i^{k_4}(p^{k_4}) &\neq g_i^{k_4}(e^{k_4}) g_i^{k_3}(p^{k_3}), \end{aligned} \tag{22a}$$

where $g_i^k(x) = s_i^k \sum_{j \in \mathcal{N}_i} a_{ij} x_j$ which only uses the entries x_j for which $j \in \mathcal{N}_i$.

Proof: Using (18)-(20), we can write Q_i as follows

$$Q_i = \underbrace{\begin{bmatrix} I & -I & \mathbf{0}_{T\times T} \\ \mathbf{0}_{T\times T} & I & -I \\ \mathbf{0}_{T\times T} & \mathbf{0}_{T\times T} & I \end{bmatrix}}_{D_i} \underbrace{\begin{bmatrix} \mathbf{a}_i & \mathbf{b}_i & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{c}_i & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{d}_i \end{bmatrix}}_{\tilde{Q}_i}.$$

Since T>1, $\tilde{\Phi}_i=\begin{bmatrix}\mathbf{a}_i & \mathbf{b}_i\end{bmatrix}$ has at least two rows, and given that (22b) holds, $\tilde{\Phi}_i$ has column rank equal to two. Moreover, if (22a) holds, \mathbf{c} and \mathbf{d} each have at least one element that is nonzero. Thus, \tilde{Q}_i has full column rank. Clearly D_i has full rank which implies that the rank of Q_i is equal to the rank of \tilde{Q}_i [11]. Therefore, there exists a unique solution to (15) using the pseudoinverse.

If one of the assumptions in (22a)-(22b) is not met, Q_i will have a nontrivial nullspace. Therefore, in that case, (21) does not have a unique solution.

Remark 1: When using the full transportation model, namely $|\mathcal{L}| > 1$, we can expand Φ_i in (18), adding two columns to Q_i for each transportation network l with entries $\sum_{j \in \mathcal{N}_i} \check{a}_{ij}^l e_j^k$ and $\sum_{j \in \mathcal{N}_i^l} \check{a}_{ij}^l p_j^k$ and the corresponding entries $\check{\beta}_i^{e,l}$ and $\check{\beta}_i^{p,l}$ to the vector on the RHS of (21). Furthermore, by similar process as shown in Corollary 1, we can construct $Q_i^{\mathcal{L}} = D_i \tilde{Q}_i^{\mathcal{L}}$ where the first row of $\tilde{Q}_i^{\mathcal{L}}$ becomes $\begin{bmatrix} \mathbf{a}_i^l, \mathbf{b}_i^l, \dots, \mathbf{a}_i^{|\mathcal{L}|}, \mathbf{b}_i^{|\mathcal{L}|}, \mathbf{0}, \mathbf{0} \end{bmatrix}$. By satisfying (22a) and ensuring that $\begin{bmatrix} \mathbf{a}_i^l, \mathbf{b}_i^l, \dots, \mathbf{a}_i^{|\mathcal{L}|}, \mathbf{b}_i^{|\mathcal{L}|} \end{bmatrix}$ has full column rank, we can also show a unique solution exists in this case.

The results in Theorem 2 and Corollary 1 allow us to learn the spreading parameters from data for homogeneous and heterogeneous viruses, respectively, under the given assumptions. Bridging these two, we can group different nodes into sets with homogeneous parameters, for example rural vs. urban counties.

V. SIMULATIONS AND CASE STUDY

In this section we apply the networked SEIR model to the COVID-19 pandemic in the Northeast US, and incorporate flight mobility data via simulations and real spread data.

A. Study area

We consider the spread of COVID-19 through five states in the Northeastern US from March through August, 2020, and consider how the underlying air transportation network between the cities in the five-state region propagated the virus. Specifically, we obtain data for New York (NY), New Jersey (NJ), Massachusetts (MA), Rhode Islands (RI), and Connecticut (CT), and consider this five-state region as a closed system (i.e., no virus entering or leaving the system). This region is selected both because this was the first significant COVID-19 outbreak in the US, making the simplifying modeling assumption that the region is closed with respect to COVID-19 more reasonable.

We model the infected population proportion in each of the 110 counties in the five-state region. Note that we combine the COVID-19 case numbers for the five counties that make up New York City into one administrative region since diagnosis statistics are provided at the city-level for New York City.

In order to capture the transmission of COVID-19 accurately, we categorize counties as either *urban* or *rural* based on the average population density. Counties with a population density of at least 500 people per square mile are considered urban, while counties with lower density are considered rural. County-level population counts are obtained from the US Census Bureau 2018 population estimate [12].

B. Transportation data and network topologies

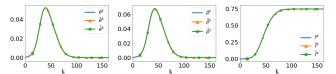
Three different types of connections are considered when modeling the network topology in (2): (i) county adjacency (a_{ij}^N) ; (ii) self-loops for spread within the county $(a_{ij}^S=I)$; (iii) flights between airports $(a_{ij}^{F,k})$ to capture long-range links between non-adjacent counties [13]. It should be noted that the flight adjacency matrix $a_{ij}^{F,k}$ is time-varying.

We incorporate travel by collecting flight data for every flight between cities in the study area from March through August, 2020 (most recent data available at time of writing) from the Bureau of Transportation (BTS) [13]. This data includes aircraft registration, which is cross-referenced with the Federal Aviation Administration database to obtain the number of passenger seats on each flight. The full dataset and code to reproduce the results are available online [14].

To construct the scaled flight matrix, we use the number of available seats between each city pair on a particular day, and normalize this number with the maximum number of daily seats observed for each city pair. This produces a scaled value between 0 and 1 that represents the intensity of travel between any given city pair on a particular day.

Phase	Start	End	ϕ_N^u	ϕ_S^u	ϕ_N^r	ϕ_S^r
1	Jan 22	Mar 26	1.0	1.0	1.0	1.0
2	Mar 27	Apr 20	0.6	0.8	1.0	1.0
3	Apr 21	Jun 8	0.6	0.8	0.8	0.85
4	Jun 9	Aug 4	1.0	1.0	1.0	1.0

TABLE I: Study phases and scaling factors for adjacency between urban counties and other counties (urban or rural) ϕ_N^u , within urban counties ϕ_S^u , adjacency between rural counties and other counties (urban or rural) ϕ_N^r , and within urban counties ϕ_S^r , inspired by observed travel reductions summarized in [15] to account for change in travel activity. The magnitude of the scaling factor reflects the extent of the connection, i.e., the higher ϕ , the stronger the connection.



(a) Network average (b) Network average (c) Network average exposed state. removed state.

Fig. 1: Simulation of a homogeneous SEIR system with three networks, its measured states, and the recovered states to show that the recovered states captures the average state.

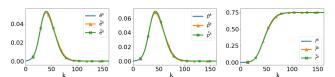
To capture the reduction in travel associated with statewide stay-at-home orders, the study period is divided into four phases shown in Table I. Phase 1 represents the time before the stay-at-home order (high transportation volume). Phase 2 represents the time immediately after the stay-at-home order started (declining transportation in urban areas). By Phase 3, both urban regions and rural regions restrict travel. Phase 4 represents a gradual return to pre-restriction travel levels. Each row of each adjacency matrix, excluding flights between airports, is scaled by the appropriate ϕ value in Table I based on its urban/rural classification and the phase to account for the corresponding reduction in mobility.

C. Simulations

In this section we illustrate the analysis and parameter estimation results from Sections III-IV.

We use the county adjacency matrix (a_{ij}^N) . To simulate the states for the SEIR model we use (1), with homogeneous spread parameters $(\beta^{e,N}, \beta^{p,N}, \sigma, \gamma) = (0.04, 0.06, 0.40, 0.30)$ and the initial state $e_1^0 = 0.02, e_2^0 = 0.03, p_1^0 = 0.01$, with the rest of the initial conditions for the non-susceptible states set to zero for each node. We correctly recover the spread parameters using (15) and e^k , p^k , and r^k for $k \in \{0,1\}$, as expected by Theorem 2.

We include the adjacency matrix that represents the flights between airports $(a_{ij}^{F,k})$, and the adjacency matrix with only self-loops for spread within the county (a_{ij}^S) . To simulate the states we use (1)-(2), with the same initial state and the spread parameters $(\beta^{e,N}, \beta^{p,N}, \beta^{e,F}, \beta^{p,F}, \beta^{e,S}, \beta^{p,S}, \sigma, \gamma) = (0.04, 0.06, 0.02, 0.03, 0.05, 0.07, 0.40, 0.30)$. Moreover, we add measurement noise to evaluate the sensitivity of the estimation results and assume that the perturbation on e is greater than that on e and e since it is the most difficult of the three states to measure. The measured states are e



(a) Network average (b) Network average (c) Network average exposed state. infected state. removed state.

Fig. 2: Simulation of a homogeneous SEIR system with three networks, one network is not completely known. Shows how well the recovered states captures the average state.

 \tilde{p} , and \tilde{r} , determined by $\tilde{e}_i^k = e_i^k + \varepsilon_e(e_i^k)$ where $\varepsilon_e(x_i) \sim$ $\mathcal{N}(0, 0.015x_i + 0.0001), \ \tilde{p}_i^k = p_i^k + \varepsilon(p_i^k), \ \text{and} \ \tilde{r}_i^k = r_i^k + \varepsilon(r_i^k)$ where $\varepsilon(x_i) \sim \mathcal{N}(0, 0.008x_i + 0.00001)$. In order to emulate the difficulty of measuring the states at the beginning of an outbreak, we start measuring from k = 14, and recover the spread parameters by left multiplying (15) by the pseudoinverse of Q. The estimated states \hat{e} , \hat{p} , and \hat{r} are constructed using (1), the first set of measured states \tilde{e}^{14} , \tilde{p}^{14} , and \tilde{r}^{14} , and the recovered spread parameters. In Figure 1 we show how well the average states are recovered compared to the average of the actual states, e, p, and r using the measured states to recover the spread parameters. The recovered spread parameters $(\hat{\beta}^{e,N}, \hat{\beta}^{p,N}, \hat{\beta}^{e,F}, \hat{\beta}^{p,F}, \hat{\beta}^{e,S}, \hat{\beta}^{p,S}, \hat{\sigma}, \hat{\gamma})$ are (0.043, 0.058, 0.023, 0.028, 0.037, 0.082, 0.400, 0.300). The error of \hat{e} , \hat{p} , and \hat{r} are 0.016, 0.015, and 0.004, respectively, computed as $\frac{\|x-\hat{x}\|_2}{\|x\|_2}$, for the corresponding state x.

To evaluate the sensitivity of recovering the states with

To evaluate the sensitivity of recovering the states with measurement noise and while only approximately knowing the network, we use a noisy version of the adjacency matrix that represents the aviation network by adding i.i.d. zero-mean Gaussian noise with standard deviation 0.001 to every possible edge, not allowing entries to be negative nor greater than 1. The recovered spread parameters $(\hat{\beta}^{e,N}, \hat{\beta}^{p,N}, \hat{\beta}^{e,F}, \hat{\beta}^{p,F}, \hat{\beta}^{e,S}, \hat{\beta}^{p,S}, \hat{\sigma}, \hat{\gamma})$ are (0.043, 0.058, 0.023, 0.025, 0.035, 0.082, 0.400, 0.300) and the error of \hat{e} , \hat{p} , and \hat{r} are 0.078, 0.074, and 0.018, respectively. In Figure 2, we see that the averages of the recovered states are fairly close to the averages of the actual states even when accurate flight data is not available.

D. Real COVID-19 spread data

We use daily COVID-19 case numbers aggregated by Johns Hopkins University (JHU) [16]. Using this dataset, we are able to estimate e_i^k , p_i^k , and r_i^k in (1). The per-capita infection rate in county i on day k, p_i^k is estimated by the number of confirmed cases in county i minus the cases that have been removed on day k and divided by the population in county i. Due to incompleteness and inaccuracies in the county-level recovery data, we estimate the state r_i^k by assuming $r_i^k - r_i^{k-1} = (p_i^{k-d_r} + r_i^{k-d_r}) - (p_i^{k-d_r-1} + r_i^{k-d_r-1})$. That is, we assume each confirmed case becomes removed after d_r days. Based on [17], we use the median recovery time $d_r = 21$ when computing the states. Due to uncertainty in the true d_r , we learn $\gamma = 1/d_r$ when calibrating the model. Similarly, we estimate e_i^k as $e_i^k - e^{k-1} = (p_i^{k+d_e} + r_i^{k+d_e}) - (p_i^{k+d_e-1} + r_i^{k+d_e-1})$. That is, the number of new exposed cases on day k equals the number of new cases that were confirmed on day

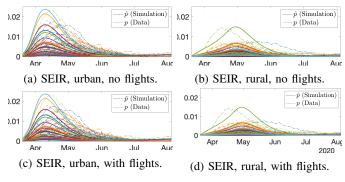


Fig. 3: SEIR model simulations both with and without taking flight data into account. Each curve represents the proportion of infected population p_i^k in a particular county. Dashed lines represent real data, p, solid lines of the same color represent the corresponding simulation results, \hat{p} .

	No flights				With flights					
	Urban		Rural		Urban			Rural		
	N	S	N	S	N	S	F	N	S	F
$\hat{\beta}^e$	5.0E-8	0.106	1.1E-8	0.124	1.4E-6	0.106	0.269	2.6E-7	0.124	7.8E-4
\hat{eta}^p	1.0E-7	0.824	1.5E-8	1.6E-7	2.3E-6	0.822	0.118	3.3E-7	4.2E-6	0.216
$\hat{\gamma}$	0.115		0.124		0.115		0.124			
Error(%)	58.6		54.1		57.7		54.1			

TABLE II: The recovered parameters and the prediction scaled error $||p-\hat{p}||_2/||p||_2$ in the case of with flight adjacency matrix and without flight adjacency matrix. E is the scientific notation.

 $k+d_e$, where d_e is the delay in number of days from becoming exposed to being confirmed. We use $d_e=14$, since COVID-19 symptoms may appear as long as 14 days after exposure [18]. Note that we use the upper bound (14 days) to account for the additional time after first showing symptoms until receiving a positive test result. To limit the number of parameters learned, we assume a fixed transition rate from exposed to infected of $\sigma=1/d_e$.

E. SEIR model with and without aviation network

Using (1)-(2) and a modified version of (21) as described in Remark 1, we estimate the parameter values and simulate the SEIR model both with and without taking transmissions resulting from inter-city travel into account. The parameters are estimated by minimizing the error in the modified version of (21) while constraining them to be non-negative using the cvx solver [19].

The SEIR model error is presented in Table II and the corresponding model performance is plotted in Figures 3a-3d. Comparing the performance of the SEIR model both with (Figs. 3c, 3d) and without (Figs. 3a, 3b) the flight network, we see that by including the aviation data, we are able to predict the proportion of the population in the infected state with slightly less error than when flight data is not considered. This indicates that, by including the transportation network, we are able to better model the virus spread. As before, the error in rural counties remains lower than in urban counties. This result is in line with our expectations that there may be viral spread over the aviation network. Note though that asymptomatic transmission is not being explicitly modeled, and may be a significant source of error in this modeling effort. Further, the inference of the epidemics states from observed data could

also be improved. Another factor that may be contributing to the higher error is assuming that the system is closed (i.e., no travel in-to or out-of the region).

VI. CONCLUSION

In conclusion, we have proposed a discrete time SEIR model to capture virus spread over transportation networks. We analyzed the limiting behavior of the model and presented conditions for estimating the spread parameters from data. The developed model is applied to infection and travel data collected from the Northeastern US. To extend this work and improve the performance of the model, we plan to incorporate asymptomatic transmission. Nonlinear state estimation could also be employed to more accurately estimate the epidemic states for the SEIR model from observed data, similar to the algorithm proposed in [10]. Capturing the transportation networks via population flows for the SEIR model, similar to [8] is another interesting future direction.

REFERENCES

- J. A. Lewnard and N. C. Lo, "Scientific and ethical basis for socialdistancing interventions against COVID-19," *The Lancet Infectious Diseases*, 2020.
- [2] M. Li and J. Muldowney, "Global stability for the SEIRS model in epidemiology," *Math. Biosciences*, vol. 125, no. 2, pp. 155–164, 1995.
- [3] X. Zhou and J. Cui, "Analysis of stability and bifurcation for an SEIR epidemic model with saturated recovery rate," *Comm. in Nonlinear Science and Numerical Sim.*, vol. 16, no. 11, pp. 4438–4450, 2011.
- [4] F. Brauer, C. Castillo-Chavez, and Z. Feng, Mathematical Models in Epidemiology. Springer, 2019.
- [5] C. Nowzari, V. M. Preciado, and G. J. Pappas, "Optimal resource allocation for control of networked epidemic models," *IEEE Transactions* on Control of Network Systems, vol. 4, no. 2, pp. 159–169, 2017.
- [6] C. Groendyke and A. Combs, "Modifying the network-based stochastic SEIR model to account for quarantine," arXiv preprint arXiv:2008.01202, 2020.
- [7] J. P. Arcede, R. L. Caga-anan, C. Q. Mentuda, and Y. Mammeri, "Accounting for symptomatic and asymptomatic in a SEIR-type model of COVID-19," arXiv preprint arXiv:2004.01805, 2020.
- [8] M. Ye, J. Liu, C. Cenedese, Z. Sun, and M. Cao, "A network SIS metapopulation model with transportation flow," in *Proceedings of the IFAC World Congress*, 2020.
- [9] R. S. Varga, Matrix Iterative Analysis. Springer-Verlag, 2000.
- [10] A. R. Hota, J. Godbole, P. Bhariya, and P. E. Paré, "A closed-loop framework for inference, prediction and control of SIR epidemics on networks," arXiv preprint arXiv:2006.16185, 2020.
- [11] R. A. Horn and C. R. Johnson, Matrix Analysis. Cambridge University Press, 2012.
- [12] US Census Bureau, "2018 population estimates by age, sex, race and hispanic origin," https://www.faa.gov/airports/planning_capacity/ passenger_allcargo_stats/passenger/media/cy19-commercial-serviceenplanements.pdf, June 2019.
- [13] Bureau of Transportation Statistics (BTS), "Bureau of transportation statistics," https://www.bts.gov/, 2020, accessed: July 2020.
- [14] M. Shang, J. Pham, D. Vrabac, B. Butler, P. Paré, and R. Stern, "Air travel data during the COVID-19 pandemic in the United States," http://hdl.handle.net/11299/217208, November 2020.
- [15] J. De Vos, "The effect of COVID-19 and subsequent social distancing on travel behavior," *Transportation Research Interdisciplinary Perspectives*, p. 100121, 2020.
- [16] E. Dong, H. Du, and L. Gardner, "An interactive web-based dashboard to track COVID-19 in real time," *The Lancet Infectious Diseases*, vol. 20, no. 5, pp. 533–534, 2020.
- [17] Q. Bi, Y. Wu, S. Mei, C. Ye, X. Zou, Z. Zhang, X. Liu, L. Wei, S. A. Truelove *et al.*, "Epidemiology and transmission of COVID-19 in 391 cases and 1286 of their close contacts in Shenzhen, China: a retrospective cohort study," *The Lancet Infectious Diseases*, 2020.
- [18] Centers for Disease Control and Prevention (CDC), "When to quarantine," https://www.cdc.gov/coronavirus/2019-ncov/if-you-are-sick/quarantine.html, 2020, accessed: Aug 2020.
- [19] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.1," http://cvxr.com/cvx, Mar. 2014.