

Latent Variables and the Learnability of Neural Networks

Ryan Bauer, Nicholas Cich, and Asa
Harbin

Introduction

- Using neural networks to learn the parameters of an ordinary differential equation system is quite difficult
 - Small amounts of noise can lead to large errors
- In an epidemiological setting, we often do not have data for impactful latent variables, which can further amplify this issue
 - Difficult to directly quantify the “S” population in an SIR model

Motivation

- If we know what latent variables are most important for determining disease spread, we can more effectively focus our data collection efforts
 - Directly collect data or indirectly collect data through proxies
- In our project, we explore how different unknown ODE system states influence the learnability of neural networks

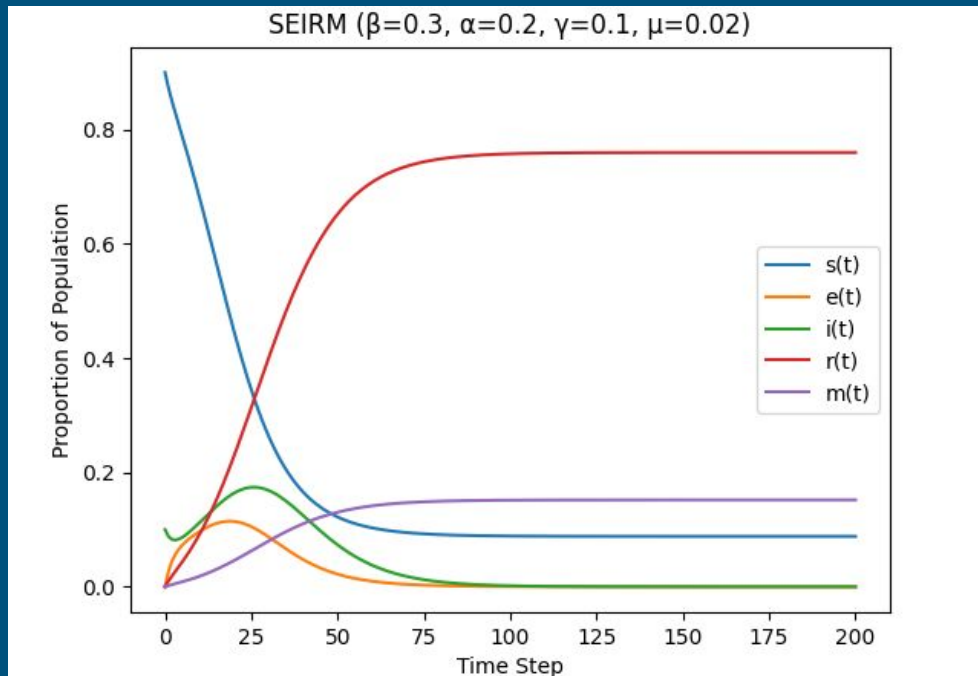
SEIRM Model

- In this project, we focused on the SEIRM compartmental model for our analysis
 - Susceptible, Exposed, Infected, Recovered, and Mortality
- SEIRM is particularly useful for analyzing diseases with long incubation periods, like COVID-19
 - Data about exposed individuals is much easier to collect before they become infected
- SEIRM is complex enough for us to see the effects of missing compartmental data, yet simple enough for us to understand quite intuitively

SEIRM Model

$$\begin{aligned}\frac{dS_t}{dt} &= -\beta_t \frac{S_t I_t}{N} & \frac{dE_t}{dt} &= \beta_t \frac{S_t I_t}{N} - \alpha_t E_t \\ \frac{dI_t}{dt} &= \alpha_t E_t - \gamma_t I_t - \mu_t I_t & \frac{dR_t}{dt} &= \gamma_t I_t & \frac{dM_t}{dt} &= \mu_t I_t\end{aligned}$$

SEIRM Model



EINNs

- Epidemiologically-Informed Neural Networks [1] are a tool designed to leverage the benefits of both mechanistic and AI-based models
 - Mathematical backing of mechanistic models
 - Ability of AI models to ingest heterogeneous information
- Data generated using epidemic mechanistic models can be used to supervise the training of neural networks
 - Networks can learn hidden dynamics of epidemic while maintaining a general framework
 - ODE equations are not numerically solved during training, but rather directly estimated

Problem Formulation

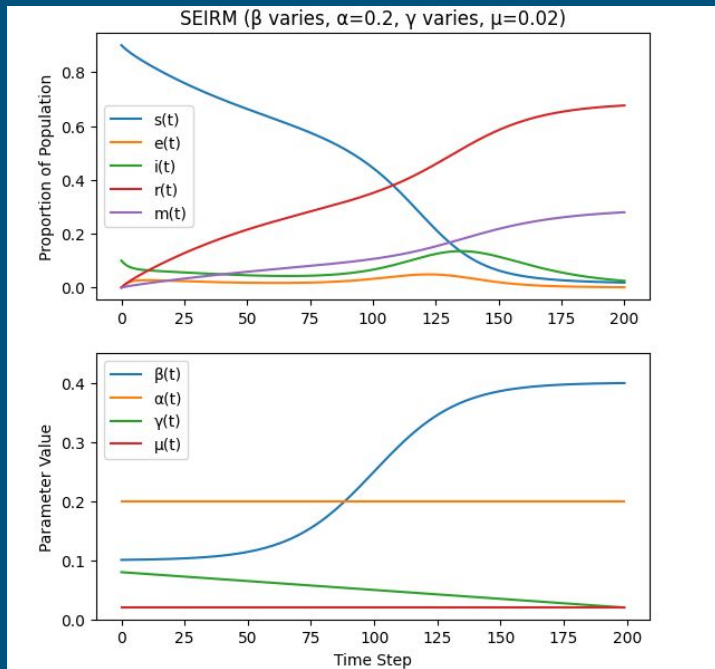
- In this project, we build upon the EINN work by Rodríguez et. al. [1] to study what data is most important when training neural networks for epidemiological models
 - Which data is the most critical for generating a high-quality estimation of dynamics of the SEIRM model?
- Can we forecast the trajectory of a disease by relying only on “infected” and “recovered” population data?

Data Generation

- For our training data, we used the SEIRM model to generate synthetic time-series data for multiple parameter combinations and record the S, E, I, R, and M values
 - Because model parameters can change as a function of time, SEIRM curves can become quite complex
 - For simplicity, we allow parameter time series to be constant, follow a linear function, or follow a logistic function
- For each parameter combination, we have five time series of length 200

$$\begin{aligned}\frac{dS_t}{dt} &= -\beta_t \frac{S_t I_t}{N} & \frac{dE_t}{dt} &= \beta_t \frac{S_t I_t}{N} - \alpha_t E_t \\ \frac{dI_t}{dt} &= \alpha_t E_t - \gamma_t I_t - \mu_t I_t & \frac{dR_t}{dt} &= \gamma_t I_t & \frac{dM_t}{dt} &= \mu_t I_t\end{aligned}$$

Data Generation Example



Experimental Setup

- First, we use an EINN to generate predictions of the S data for an SEIRM model based on the E, I, R, and M time series
 - We then compare the predicted time series to the actual time series
- We repeat this process for each of the remaining time series, “covering” the data from one compartment at a time
- Lastly, we compare the SSEs of each predicted time series to determine the importance of each compartment

Results

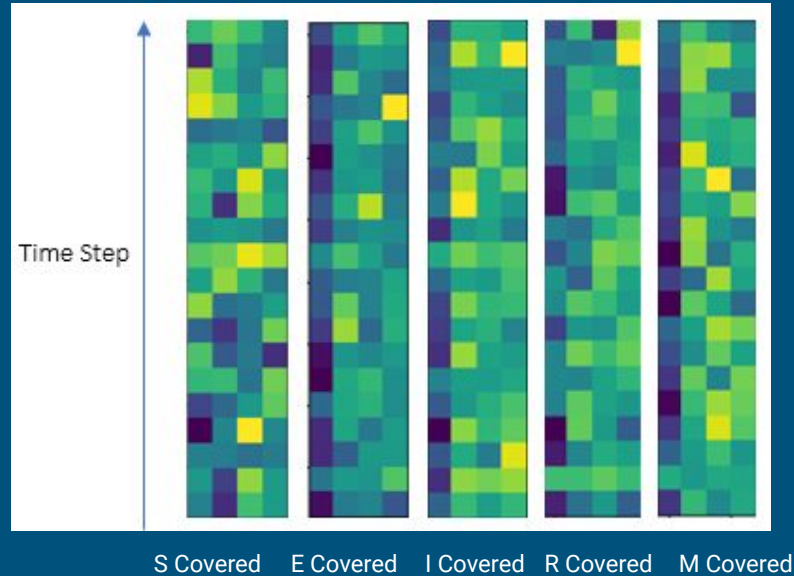
Curve Covered	Average Loss (25 Randomly Generated Datasets)
S(t)	0.0057
E(t)	0.0029
I(t)	0.0091
R(t)	0.0106
M(t)	0.0064

Model Explainability

- In addition to learnability, we also explore model explainability through integrated gradients
 - Given a trained neural network and an input, we find the most relevant input features by summing the difference between the gradient of the input and the gradient of informationless baseline
- This technique allows us to qualitatively analyze the importance of various compartments in the SEIRM model

$$\text{IntegratedGrads}_i(x) := (x_i - x'_i) \times \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_i} d\alpha$$

Model Explainability Results



Conclusions

- Susceptible population data does not seem as crucial to collect as we initially anticipated
- Infected and recovered population data is significantly more important than the remaining compartmental data
- Exposed population data seems to be the easiest to predict given the remaining four curves

Future Work

- Test how different levels of noise factor into the importance of each group
 - Could have implications for priorities based on how reliable the data collections may be in a real world scenario.
- Try combinations of covered variables
- Test out results on a real world data set instead of synthetic data
- Test experiment with different types of Epidemiology focused Neural Networks

References

- [1] Rodríguez, A., Cui, J., Ramakrishnan, N., Adhikari, B., and Prakash, B. A. (2022). EINNs: Epidemiologically-Informed Neural Networks. arXiv preprint arXiv:2202.10446.