# Interpreting County Level COVID-19 Infection and Feature Sensitivity using Deep Learning Time Series Models

Md Khairul Islam [1], Di Zhu [1], Yingzheng Liu [1], Andrej Erkelens [2] , Nick Daniello [2], Judy Fox [1,2]

[1] Computer Science Department, University of Virginia

[2] School of Data Science, University of Virginia

Charlottesville, USA

Email : {mi3se, yqx8es, yl4dt, wsw3fa, njd9e, cwk9mp}@virginia.edu

*Abstract*—Interpretable machine learning plays a key role in healthcare because it is challenging in understanding feature importance in deep learning model predictions. We propose a novel framework that uses deep learning to study feature sensitivity for model predictions. This work combines sensitivity analysis with heterogeneous time-series deep learning model prediction, which corresponds to the interpretations of Spatio-temporal features from what the model has actually learned. We forecast county-level COVID-19 infection using the Temporal Fusion Transformer (TFT). We then use the sensitivity analysis extending Morris Method to see how sensitive the outputs are with respect to perturbation to our static and dynamic input features. The significance of the work is grounded in a real-world COVID-19 infection prediction with highly non-stationary, finely granular, and heterogeneous data. 1) Our model can capture the detailed daily changes of temporal and spatial model behaviors and achieves high prediction performance compared to a PyTorch baseline. 2) By analyzing the Morris sensitivity indices and attention patterns, we decipher the meaning of feature importance with observational population and dynamic model changes. 3) We have collected 2.5 years of socioeconomic and health features over 3142 US counties, such as observed cases and deaths, and a number of static (age distribution, health disparity, and industry) and dynamic features (vaccination, disease spread, transmissible cases, and social distancing). Using the proposed framework, we conduct extensive experiments and show our model can learn complex interactions and perform predictions for daily infection at the county level. Being able to model the disease infection with a hybrid prediction and description accuracy measurement with Morris index at the county level is a central idea that sheds light on individual feature interpretation via sensitivity analysis.

*Index Terms*—Interpretability, County Level COVID-19, Time Series Deep Learning, TFT, Sensitivity Analysis, Morris Method.

## I. INTRODUCTION

Interpretation of machine learning models has recently [1] led to numerous research applications of AI for social impact. This includes direct analysis of model components with casual inference and uncertainty estimation or studying sensitivity to input perturbations. Typically a simpler model is easier to interpret but can result in lower predictive accuracy. One natural question that arises is how to interpret these complex deep learning models, which may describe the data better. One major challenge of interpretability is the gap between model prediction accuracy and descriptive accuracy in real-world problems. The latter can be illustrated by a quantifiable measurement and explanation of the individual feature importance with regard to the model's forecast relevancy.

To our knowledge, however, no prior studies have evaluated individual feature importance at the county level using deep learning and the Morris method. We have been closely monitoring the scientific literature and identifying reports describing the community-level impact of COVID-19. A number of factors contribute to COVID-19 cases and deaths, including a very diverse set of socioeconomic and geographic-specific features. A more granular real-time analysis that considers important county-level factors is lacking and urgently needed. Furthermore, non-stationary time series (with their distribution drifting over time) [2] or time series with extreme events [3] or unknown events like COVID variants are particularly challenging to model and interpret.

To effectively study county-level input features, we design a novel method to compute the Morris index but generalize it to multidimensional spatial and temporal variables. Using a self-attention-based Temporal Fusion Transformer (TFT) model [4], we can capture a complex mix and full range of static and dynamic covariates, known inputs, and other exogenous time series parameters. We perform individual feature importance evaluations to identify the most influential features for prediction and the sensitivity of infected cases. The results show that the model obtains significant performance and learns temporal patterns. More significantly, our scaled Morris index provides sensitivity measurement to individual features that help policymakers develop effective control strategies in response to the rapidly evolving pandemic. We have made our code available on GitHub [1]. In summary, we've made the following contributions:

- Introduce individual feature sensitivity to forecasting outputs with an extended Morris Method for multidimensional spatial and temporal data.

[1] https://github.com/Data-ScienceHub/gpce-sensitivity

- Model heterogeneous time-series prediction and analyze attention weights for insights on feature importance.
- Stratify county-level population characteristics(Age and Industry segments) from socioeconomic and health data.

The rest of the paper is structured as follows. Section II discusses the data collection, and feature descriptions. Section III presents the background on the TFT model architecture and the Morris method. Section IV describes the data pre-processing and experimental setups. Section V analyzes the temporal patterns and feature importance insight from the TFT. Section VI discusses the sensitivity analysis with the Morris method. Then Section VII discusses the related works and Section VIII has the concluding remarks and impact on possible future works.

## II. INPUT DATA AND FEATURES

We collected our dataset for 3142 US counties. They are from multiple sources, including CDC (Centers for Disease Control and Prevention), USA Facts [5], Unacast [6]. The dynamic features include entries from 02-29-2020 to 05-17-2022. Except for vaccination, where the earliest available data in CDC [7] was from 12-14-2020 when the US initiated a nationwide COVID-19 vaccination campaign. In total we select 9 observed features, static and dynamic, to predict cases and deaths. Fig.1 summarizes the feature groups with the influencing factors they capture and which county characteristics they represent.
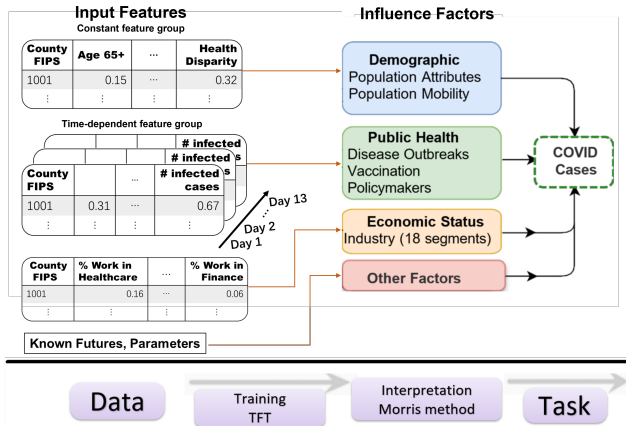


Fig. 1: The data: the feature groups and influencing factors.

Table I lists the features used by our model and the respective sources with descriptions. In particular, age distribution, health disparities, disease spread, social distancing, and transmissible cases features are collected from the outputs of the COVID-19 Pandemic Vulnerability Index (PVI) dashboard [8], maintained by the National Institute of Health (NIH).

## III. BACKGROUND AND THEORETICAL FOUNDATION

### A. Temporal Fusion Transformer

We used the TFT model [4] to predict daily COVID-19 cases and deaths at the county level. For this work, we dive deeper into the COVID-19 daily cases prediction and

combine sensitivity analysis of individual features. Figure 2 shows a high-level overview of the work. Gated Residual Network (GRN) is the building block of TFT and it enables more efficient use of the model architecture. TFT takes static metadata, time-varying past inputs, and time-varying known future inputs. The model inputs are passed through a Variable Selection Network (VSN) to select the most salient features and filter out noise.
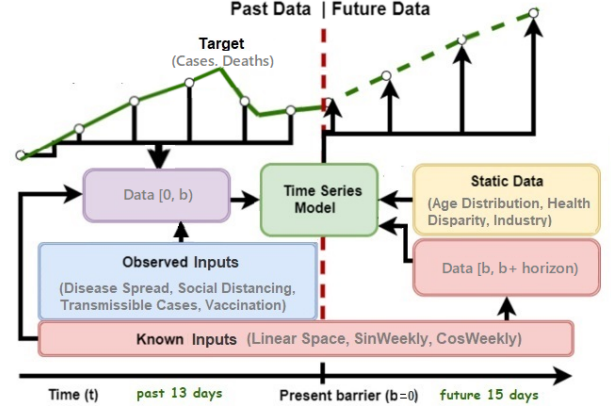


Fig. 2: A time series forecasting model. Each sliding window consists of time-sequential data that is split into two parts, the past, and the future.

Learning significant data points is done by leveraging local context with LSTM-based sequence to sequence layer. Past inputs are fed into the encoder, whereas known future inputs are fed into the decoder. Their outputs go through a static enrichment layer which enhances temporal features with static metadata. Following static enrichment, TFT adds a novel interpretable multi-head self-attention mechanism to better learn the different temporal patterns. This allows TFT to learn long-rage dependencies that can be challenging for Recurrent Neural Network (RNN) based models. Following the self-attention layer, additional gating layers are added to facilitate training.

### B. Sensitivity Analysis and The Morris Method

Sensitivity Analysis is the study of the input-output relationship in a computational model [11]. It can identify the importance of each model parameter in determining the outputs. [12] proposed gradient-based attribution which approximates the neural network function $f(\boldsymbol{X})$ around a given input $\boldsymbol{X}$ by linear part of the Taylor expansion as

$$f(\boldsymbol{X} + \Delta \boldsymbol{X}) \approx f(\boldsymbol{X}) + \nabla_{\boldsymbol{X}} f(\boldsymbol{X})^T \cdot \Delta \boldsymbol{X} \qquad (1)$$

and they analyzed the network sensitivity by looking at how small changes $\Delta \boldsymbol{X}$ at the input correlate with changes at the output. Gradient $\nabla_{x_i} f = \frac{\delta f(\boldsymbol{X})}{\delta}$ gives the linear approximation of this change for a change in the $i$-th input token $x_i \in \mathbf{R}$, and the attribution of how much input token $x_i$ affects the network output $f(\boldsymbol{X})$ can be approximated by the L2 norm of $\nabla_{x_i} f$.

TABLE I: Short description of input feature groups (and target features). Refer to references and Appendix for full details.

| Feature | Description | Data Source | Input Type |
|---|---|---|---|
| Cases | Daily COVID-19 cases | USA Facts [5] | Target |
| Deaths | Daily COVID-19 deaths | | |
| Age Distribution | Percentage of population aged 65 or older | SVI [9] | Static |
| Health Disparities | Uninsured population percent and socioeconomic status | | |
| Industry | Percentage of population in different industry sectors (only used in Section VI) | Census Bureau [10] | |
| Vaccination | Percentage of population fully vaccinated | CDC [7] | Observed |
| Disease Spread | Fraction of total cases from the last 14 days (one incubation period) | USA Facts [5] | |
| Transmissible Cases | Population size divided by cases from the last 14 days | USA Facts [5] | |
| Social Distancing | Change in distance travelled relative to baseline(previous year), based on cell phone mobility data | Unacast [6] | |
| SinWeekly | sin (day of the week/7) | Date | Known Future |
| CosWeekly | cos (day of the week/7) | Date | |
| Linear Space | Unique index for each county. | USA Facts [5] | |

## C. Problem Statement

We will use deep learning to study feature sensitivity for model predictions of COVID-19 infection at the county level. Given this, we adopted the Morris method [13], a reliable and efficient sensitivity analysis method that defines the sensitivity of a model input as the ratio of the change in an output variable to the change in an input feature. More precisely, given a model $Y = f(X)$, the sensitivity (or the elementary effect) of a model input feature $x_i$ can be defined as

$$EE_i(X) = \frac{y(x_1, x_2, \ldots, x_i + \Delta, \ldots, x_k) - y(X)}{\Delta} \quad (2)$$

where $X$ is a scaled vector of $k$ parameters and $\Delta$ is the change to an input feature. Since elementary effects may cancel each other out, the mean of the absolute values in distribution $EE_i(X)$, denoted by $\mu^*$ (called **the Morris Index**), is recommended because it provides true importance of features [14].

---

**Algorithm 1:** Novel Morris Index Calculation for **Spatio-temporal** data

---

**Input:** $X = \{x_1, x_2, \ldots, x_k\}$, target feature $x_i \in X$ with dimension $[C, T]$, model $y$, $\Delta$
// $X$ is a set of $k$ input features, $\Delta$ is the change to $x_i$

1   $Y_\Delta = y(x_1, x_2, \ldots, x_i + \Delta, \ldots, x_k)$
2   $Y = y(X)$
3   **while** $t < T$ **do** /* Temporal */
     // Loop through 640 Days
4     **while** $c < C$ **do** /* Spatial */
       // Loop through 3142 US Counties
5       $G \leftarrow G + |Y_\Delta[c][t] - Y[c][t]|$ /*Total Change*/
6       $c \leftarrow c + 1$
7     $t \leftarrow t + 1$
8     $c \leftarrow 0$
   // Calculate normalized Morris Index $\hat{\mu}^*$
9   $\hat{\mu}^* = G/(C * T * \Delta)$
10   **return** $\hat{\mu}^*$

---

The original Morris method was proposed to screen static input factors (or static features) but this is not conventional for the time series dataset (or dynamic features) where we look at time variation and spatial variation with an important overall influence on the output of COVID-19 cases prediction using TFT. Hence, we design and implement a revised Morris Algorithm 1 to handle the Spatio-temporal COVID-19 data sequences. The algorithm calculates a **normalized Morris Index** $\hat{\mu}^*$ by dividing the total change to the output $G$ by the total number of counties $C$, the total number of daily timestamps $T$ and the change to the input $\Delta$. In this study, $C$ is the total number of counties and $T$ is the total number of daily timestamps between 2-29-2020 and 11-29-2021.

## IV. EXPERIMENTAL SETUP

### A. Computational Resources

We implement our TFT model with both Tensorflow [4] and PyTorch [15]. Then we conducted a performance evaluation of the model training on Google Colab and HPC clusters including the GPU nodes in Table II. The model training time is about 30 hours. Each training epoch takes on average 50 minutes on a GPU node with at least 32GB of RAM. Each Morris runs with a trained model, and with additional feature analysis that takes around 35 minutes.

TABLE II: Runtime environment and hardware specification.

| Driver | CUDA | Processor | NVIDIA GPU |
|---|---|---|---|
| 470.82.01 | 11.4 | Intel Xeon | A100-SXM4-40GB |
| | | | Tesla P100-PCIE |
| | | | Tesla V100-SXM2 |
| | | | Tesla K80 |

### B. Evaluation Metrics

Our forecasting models are evaluated using the following metrics. Mean Squared Error (MSE) is used as the loss function following prior works on COVID-19 forecasting [2], [16]. Other metrics include Mean Absolute Error (MAE), Root Mean Square Error (RMSE), Symmetric Mean Absolute Percentage Error (SMAPE), and Normalized Nash-Sutcliffe Efficiency (NNSE) [17].

These metrics have been widely used in evaluating regression model performance [18]. The benefit of using NNSE is its robustness to error variance. NNSE is 1 for a perfect model. A model with an error variance equal to that of the observed time series will give NNSE = 0.5 (NSE=0). When the error

variance is larger, NNSE will be in the range (0, 0.5). Also note that MAE favors small population counties more than MSE, where large population counties dominate the error evaluation. SMAPE is very sensitive to the quality of model fitting in small counties. Hence, we can order the metrics (RMSE, NNSE, MAE, and SMAPE) in their evaluation reliability.

## C. Pre-processing

The original data contains outliers caused by rare events or human errors in data collection. We clean the outliers from our input features using the following lower and upper thresholds,

$$
\begin{aligned}
lower &= Q1 - (7.5 * IQR) \\
upper &= Q3 + (7.5 * IQR)
\end{aligned} \tag{3}
$$

where $Q1, Q3$ are the first and third percentiles and $IQR$ is the interquartile range. The data statistics before and after removing the outliers are reported in Table III and the ground truth is plotted in Figure 3. Both input and target features were min-max scaled before feeding to the model. We did not use the moving average to further smooth this dataset because that would filter out persistent temporal patterns observed in the raw data.

TABLE III: Statistics of input features. Cases and vaccination have much higher variance than others.

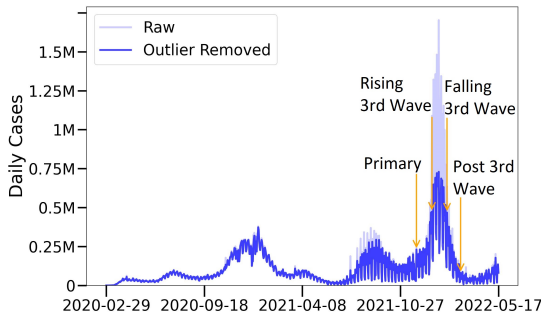| Feature | Original | | Cleaned | |
|---|---|---|---|---|
| | **Mean** | **Std** | **Mean** | **Std** |
| **Cases** | **31.67** | **337.4** | **27.18** | **174.2** |
| Deaths | 0.378 | 2.853 | 0.239 | 2.244 |
| Age Distribution | 0.576 | 0.094 | 0.576 | 0.094 |
| Health Disparities | 0.368 | 0.198 | 0.368 | 0.198 |
| **Vaccination** | **20.61** | **22.92** | **20.61** | **22.92** |
| Disease Spread | 0.150 | 0.194 | 0.150 | 0.193 |
| Social Distancing | 0.784 | 0.228 | 0.795 | 0.229 |
| Transmissible Cases | 0.492 | 0.210 | 0.491 | 0.210 |



Fig. 3: Ground truth of reported COVID-19 cases [5] along with dataset splits.

## D. Benchmark

We compare our TFT model performance with a simple PyTorch *Baseline* model that predicts future cases and deaths by repeating the last known observations. So there are no parameters to learn for this baseline model. Note that the *std* of "Cases" in our dataset is widely volatile with factors 10X and 5X for raw and outlier removed data respectively. This implies that prediction without 7-day averaging (as in our benchmark) will be 2X to 5X variation from the ground truth, while the smoothed data will be varying by a factor of 2X to 3X at the peak for county-level daily cases prediction.

*a) Train-Validation-Test Split:* We partitioned our dataset into train, validation, and testing sets following the split reported in Table IV. All experiments in this work follow the primary split unless explicitly mentioned. The COVID-19 cases in the US peaked on January 15, 2022, due to new COVID-19 variants [5] like delta and omicron. This period is labeled as the third COVID-19 wave [19] and we named our additional dataset splits based on the part of this third wave included in it.

TABLE IV: Dataset splits and dates.

| Split | Start | End | Dataset |
|---|---|---|---|
| | 02-29-2020 | 11-29-2021 | Train |
| Primary | 11-30-2021 | 12-14-2021 | Validation |
| | 12-15-2022 | 12-29-2022 | Test |
| | 02-29-2020 | 12-31-2021 | Train |
| Rising 3rd Wave | 01-01-2022 | 01-15-2022 | Validation |
| | 01-16-2022 | 01-30-2022 | Test |
| | 02-29-2020 | 01-31-2021 | Train |
| Falling 3rd Wave | 02-01-2022 | 02-15-2022 | Validation |
| | 02-16-2022 | 03-02-2022 | Test |
| | 02-29-2020 | 02-28-2021 | Train |
| Post 3rd Wave | 03-01-2022 | 03-15-2022 | Validation |
| | 03-16-2022 | 03-30-2022 | Test |

*b) Model Parameters and Tuning:* Our TFT models take the prior 13 days of data as input and use that to predict both cases and deaths for the next 15 days. Therefore, at least 13 days of observed data are needed to start giving predictions. Table V presents the parameters for the TFT model and a list of hyper-parameter tuning using the validation set.

TABLE V: TFT model training and network parameters. For hyper-parameter tuning, the list is added and optimal values are in bold.

| Parameters | Value | Parameters | Value |
|---|---|---|---|
| learning rate | [**1e-3**, 1e-4] | batch size | 64 |
| hidden layer size | [**16**, 32, 64] | dropout rate | 0.2 |
| attention head size | [1, **4**] | optimizer | adam |
| gradient clip norm | [0.01, **1.00**] | loss | MSE |

*c) Training and Prediction Performance:* Table VI shows that TFT outperforms the baseline model in all metrics. A lower score is better for MAE, RMSE, and SMAPE. Higher is better for NNSE. The losses are calculated using MSE at individual county and date levels. The aggregated prediction plots are presented in Figure 4.

TABLE VI: Training and test results and comparison.

| Target | Model | MAE | RMSE | SMAPE | NNSE |
|---|---|---|---|---|---|
| Cases | TFT | 37.40 | 237.0 | 0.892 | 0.649 |
| | Baseline | 48.47 | 269.0 | 0.921 | 0.588 |
| Deaths | TFT | 0.231 | 1.36 | 0.121 | 0.648 |
| | Baseline | 0.346 | 2.25 | 0.139 | 0.401 |

We further report two experiments: 1) Geo-spatial partition with representative large and small counties, 2) Temporal partition with three additional dataset splits. Figure 5 shows how the TFT model performs for individual counties with different population sizes. The small counties often have zero COVID-19 cases reported, whereas large counties show sudden large spikes. These features make the prediction more challenging.

The test results are presented in Figure 6 for the additional splits, after extending the dataset to more recent dates as reported in Table IV. This shows that the model performs consistently during different phases of COVID-19 waves.

## V. Attention Weights and Analysis

Here we present the interpretability use cases of the TFT model in predicting COVID-19 infections. We demonstrate the use cases by (1) analyzing input variable importance, and (2) learning temporal patterns from the dataset. This gives more insight into the significant features and seasonality TFT learns from the dataset.

### A. Persistent Temporal Patterns

Analyzing persistent temporal patterns is a key to understanding the time-dependent relationships present in a given dataset. To improve the interpretability of such patterns the TFT model uses modified multi-head attention that shares values in each head and additively aggregates them:

$$\text{InterpretableMultiHead}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \tilde{\boldsymbol{H}} \boldsymbol{W}_H \quad (4)$$

$$
\begin{aligned}
\tilde{\boldsymbol{H}} &= \tilde{A}(\boldsymbol{Q}, \boldsymbol{K}) \boldsymbol{V} \boldsymbol{W}_V \\
&= \left\{ 1/H \sum_{h=1}^{m_H} A\left(\boldsymbol{Q}\boldsymbol{W}_Q^{(h)}, \boldsymbol{K}\boldsymbol{W}_K^{(h)}\right) \right\} \boldsymbol{V} \boldsymbol{W}_V \\
&= 1/H \sum_{h=1}^{m_H} \text{Attention}(\boldsymbol{Q}\boldsymbol{W}_Q^{(h)}, \boldsymbol{K}\boldsymbol{W}_K^{(h)}, \boldsymbol{V}\boldsymbol{W}_V)
\end{aligned}
\quad (5)
$$

where $\boldsymbol{W}_Q^{(h)} \in \mathbf{R}^{d_{model} \times d_{attn}}$, and $\boldsymbol{W}_K^{(h)} \in \mathbf{R}^{d_{model} \times d_{attn}}$ are head-specific weights for keys ($K$) and queries ($Q$). A() is a softmax normalization function and $m_H$ is the number of attention heads. $\boldsymbol{W}_V \in \mathbf{R}^{d_{model} \times d_v}$ is for weights shared across all heads and $\boldsymbol{W}_H \in \mathbf{R}^{d_{attn} \times d_{model}}$ is for final linear mapping.

When representing the weights of our self-attention mechanism from equation 5, we are left with a $(\boldsymbol{H}, n_s, d_s, d_s)$ array, where $\boldsymbol{H}$ is the vector of attention heads, $n_s$ is the total number of sequences in our dataset, and $d_s$ is the combined input and output sequence length. $n_s$ can be calculated by taking $n_l * (d - s + 1)$, where $d$ is the number of time steps per ID, and $n_l$ is the number of IDs (which is **counties or FIPS in the COVID-19 prediction**). By taking the mean over $\boldsymbol{H}$ and $n_s$ we obtain $\bar{\boldsymbol{H}}$, which leaves us with a high-level view of temporal patterns in each sequence. Left with an $(d_s, d_s)$ array, which is a lower triangular matrix because of the attention mask, we can see the weight at each available day $d_c$ attributed to any other visible day $d_{c-i}$ in the sequence. There is an illustration of this in Figure 7.

Figure 8a presents the mean self-attention weights for the past 13 days of data on the train set for the one-step-ahead forecast. There is a clear weekly pattern, which peaks on the same day (as the present day) in the previous week (time index -7). Attention increases again as the past day gets closer to the present day (time index 0). Implying the most recent day in the past (time index -1) also has a significant impact on the prediction. Figure 8b presents the average daily self-attention weights on the train set for the one-step-ahead forecast, along with COVID-19 cases and annotated holidays. We choose the following federal holidays within this highlighted period: (1) Thanksgiving - Nov 26, (2) Christmas Eve - Dec 25, (3) New Year's Day - Jan 1. It shows that the reported cases have a weekly pattern, which often peaks on Friday and bottoms on Sunday. Similarly, during the holidays, there is a drop in reported cases. The TFT learns such seasonality patterns from the data and the attention weights also show a similar trend where it bottoms on Friday, right before the start of the weekend, and also around the holidays. Then it peaks around Wednesday as we can see a rise in reported cases on Thursday and Friday.

### B. Input Variable Importance

We analyze the variable importance by summing up the weights from the Variable Selection Network (VSN) for each variable across the train set. The weights for each feature type are normalized to percentage and presented in Table VII. For static covariates, the highest weight is given to age distribution, we know the elder population was severely affected by COVID-19. In the observed inputs, past target values (cases, deaths) are the most important as expected. In known future inputs, the county identifier (linear Space) and day of the week (sinweekly, cosweekly) have balanced importance as they all help to learn persistent spatial and temporal patterns.

TABLE VII: Feature importance from variable selection weights. The highest values are highlighted in bold.

| Feature | Static | Observed | Known |
|---|---|---|---|
| **Cases** | | **34.47%** | |
| Deaths | | 15.94% | |
| **Age Distribution** | **62.04%** | | |
| Health Disparities | 37.96% | | |
| Vaccination | | 7.90% | |
| Disease Spread | | 6.22% | |
| Transmissible Cases | | 4.49% | |
| Social Distancing | | 4.85% | |
| SinWeekly | | 13.83% | 37.04% |
| CosWeekly | | 8.66% | 23.36% |
| **Linear Space** | | 3.65% | **39.60%** |

## VI. Feature Sensitivity Analysis

### A. Model Sensitivity of Individual Feature

Our TFT model is trained using a variety of features, shown in Table I. Intuitively, these features should have different impacts on the infection of COVID-19. To detect such influence, we studied the sensitivity of individual features using the Morris method. The normalized Morris index $\hat{\mu}^*$ of each
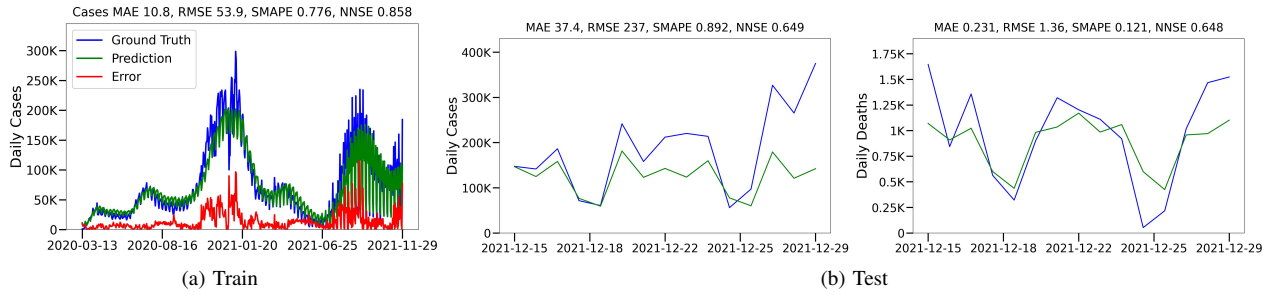
Fig. 4: TFT performance based on the primary split (aggregated over 3142 US counties). Case training performance is reported since it is used in the feature sensitivity analysis.
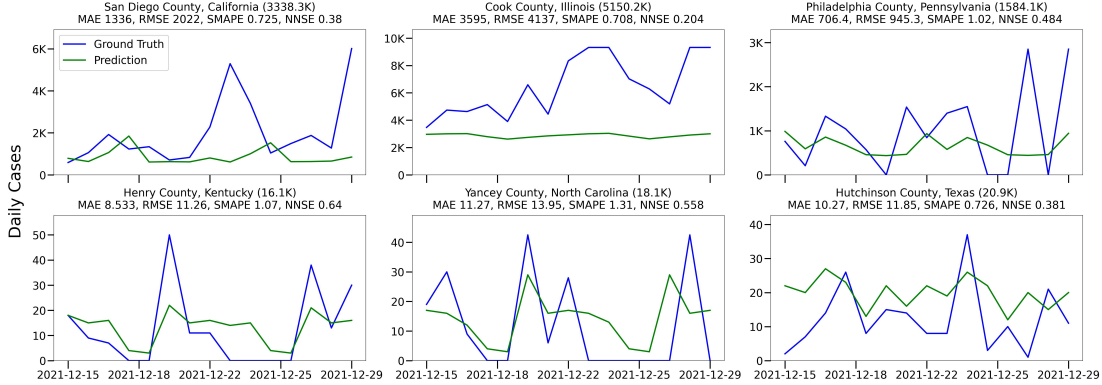


Fig. 5: Spatial: Cases prediction performance for six randomly selected US counties. The top row contains counties selected from the Top 100 US counties by population. The bottom row counties are selected from the rest. The county population is reported within brackets.
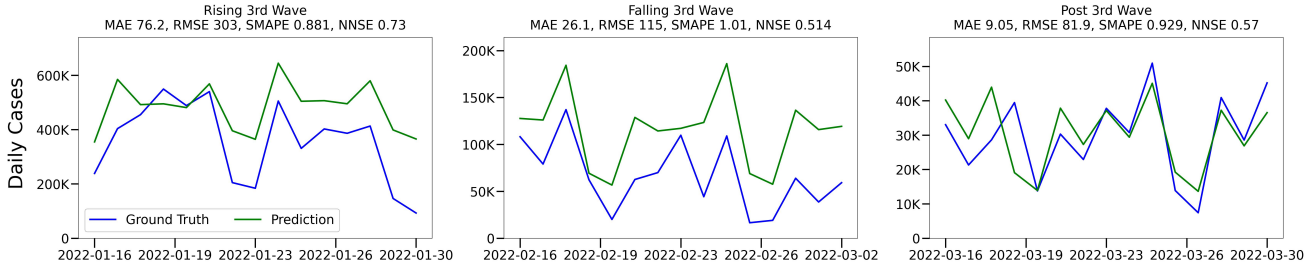


Fig. 6: Temporal: Cases prediction performance on all US counties for additional test data splits.

feature is computed which is introduced in Section III. Since the counties differ from each other drastically in the feature space, to measure feature importance in a more accurate way, we scale $\hat{\mu^*}$ by the standard deviation of each input feature $\sigma_i$, and $\hat{\mu^*} * \sigma_i$ is used to measure the sensitivity of input features on each county at a single timestamp. We refer to $\hat{\mu^*} * \sigma_i$ as **Scaled Morris Index** in the following texts. Features with a higher scaled Morris index have a greater impact on the output.

Figure 9 is for our sensitivity analysis experiments. Vaccination and disease spread are the most influential features, while other features have much less influence on the output. This result gives proof of the impact of vaccination on COVID-19.

### B. Sensitivity of Vaccination to Cases by Feature Subgroups

In the following two subsections, we discuss in detail our Morris experiments and the results on population in age subgroups and population in industry sectors as static features respectively. While both vaccination and disease spread are influential dynamic features, we choose to focus on studying the impact of vaccination on finer-grained population groups. In this study, we experimented with 2 population segmentation methods: 1) Population segmentation by age, and 2) Population segmentation by industry sectors.

A total of 7 age subgroups and 11 industry sectors were evaluated in our experiments. We trained our model for each of the population subgroups as the static input and vaccination as the dynamic input, forming a set of 18 individually trained
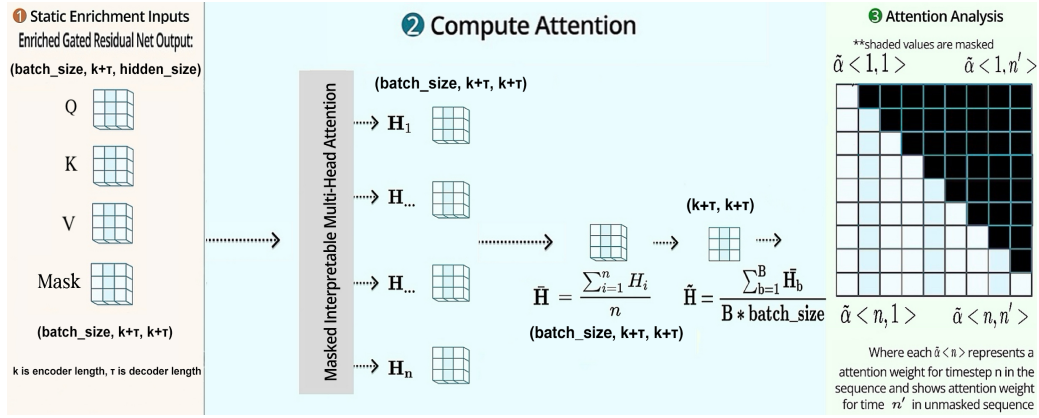
Fig. 7: Flow of aggregation and selection for TFT attention weights.



(a) Attention weights aggregated by past time index.



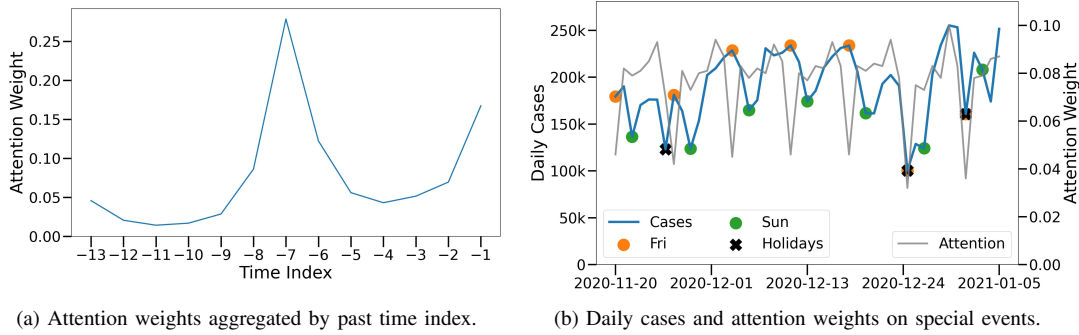(b) Daily cases and attention weights on special events.

Fig. 8: Persistent temporal patterns for one-step-ahead forecast. (a) Report from the same day in the previous week (time index -7) is the most critical in predicting COVID-19 infection. The most recent day in the past is important too, as attention increases close to the present day (time index 0). (b) The attention drops before weekends and holidays when fewer cases are reported.
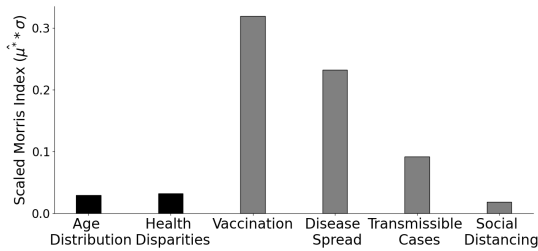


Fig. 9: Sensitivity of observed features

TABLE IX: Sensitivity analysis models and training loss

| Feature Subgroup | Training Loss | $\hat{\mu}^* * \sigma$ ($\Delta = 0.005$) |
|---|---|---|
| 0-19 | 1.50E-04 | 7.20E-07 |
| 20-29 | 1.50E-04 | 1.05E-02 |
| 30-39 | 1.50E-04 | 2.95E-03 |
| 40-49 | 1.49E-04 | 2.68E-04 |
| 50-64 | 1.51E-04 | 6.22E-05 |
| 65-79 | 1.50E-04 | 8.96E-03 |
| 80+ | 1.52E-04 | 5.97E-04 |
| Healthcare | 1.49E-04 | 3.53E-03 |
| Manufacturing | 1.61E-04 | 7.92E-03 |
| Retail | 1.50E-04 | 3.26E-03 |
| Accommodation | 1.51E-04 | 5.43E-03 |

models using TFT with the same experimental setup as described in section IV. Table IX lists the feature combinations and the training losses of a subset of the trained models.

### C. Sensitivity of Vaccination by Age Population Groups

CDC evaluated the association between county-level COVID-19 vaccine uptake and rates of COVID-19 cases and deaths in the US from January 1, 2020, to October 31, 2021 [21]. The results showed US counties with a percent of persons with $\geq$ 80% of their residents $\geq$12 years of age fully vaccinated against COVID-19 had 30% and 46% lower rates of

COVID-19 cases and death respectively. Our sensitivity study in Fig. 9 also identifies vaccination as the most influential feature. In accordance with CDC's recent evaluation [21], we divided the population into 7 age groups reported in Table VIII: 0-19, 20-29, 30-39, 40-49, 50-64, 65-79, 80+. The percentage of the population in each age group across all US counties is reported in Fig. 10.

Fig. 11 shows Morris experiments and results on the age subgroups. The solid lines specify working age groups and dotted lines represent the non-working age groups. Figure 11a

TABLE VIII: County level age statistics [20].

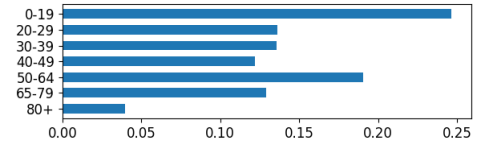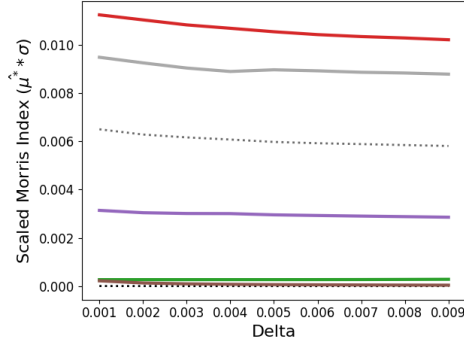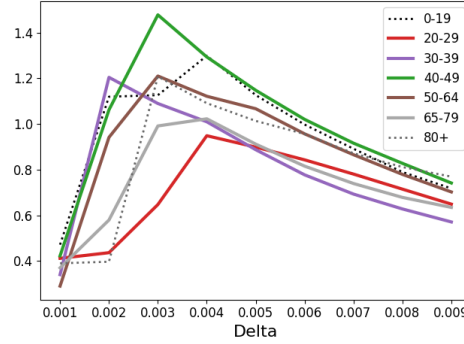| Age | Mean | Std | Ratio |
|------|-------|-------|-------|
| 0-19 | 0.242 | 0.036 | 6.622 |
| 20-29 | 0.121 | 0.030 | 3.997 |
| 30-39 | 0.118 | 0.017 | 6.854 |
| 40-49 | 0.114 | 0.012 | 8.886 |
| 50-64 | 0.200 | 0.023 | 8.633 |
| 65-79 | 0.153 | 0.036 | 4.216 |
| 80+ | 0.049 | 0.015 | 3.145 |



Fig. 10: Percent of age group in US.



(a) Age



(b) Vaccination

Fig. 11: Age subgroup sensitivity results (the greater the value, the more important the feature sensitivity). We study the features that measure a fraction of the county population by age. We then plot the scaled Morris index that multiplies the derivative with the standard deviation but only for those features with larger means.

TABLE X: County level industry statistics [22].

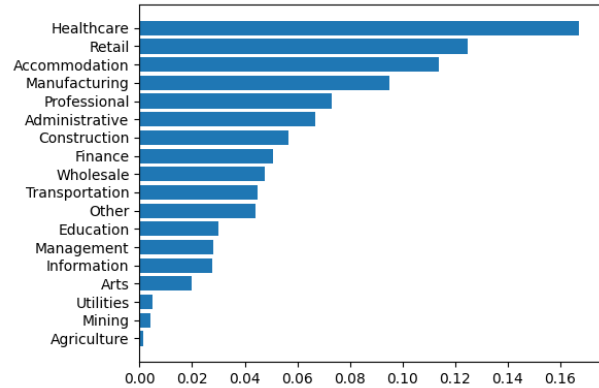| Industry | Mean | Std | Ratio |
|----------|------|-----|-------|
| **Health Care and Social Assistance** | **0.18** | **0.09** | 0.5 |
| **Manufacturing** | **0.15** | **0.13** | 0.85 |
| **Retail Trade** | **0.16** | **0.07** | 0.42 |
| **Accommodation and Food Services** | **0.12** | **0.07** | 0.62 |
| Mining, Quarrying, and Oil and Gas Extraction | 0.02 | 0.06 | 3.6 |
| Construction | 0.06 | 0.05 | 0.87 |
| Transportation and Warehousing | 0.04 | 0.05 | 1.08 |
| Professional, Scientific, and Technical Services | 0.04 | 0.04 | 1.05 |
| Wholesale Trade | 0.04 | 0.05 | 1.08 |
| Administrative and Support | 0.04 | 0.05 | 1.29 |
| Educational Services | 0.01 | 0.03 | 2.16 |
| Utilities | 0.01 | 0.03 | 3.35 |
| Finance and Insurance | 0.04 | 0.03 | 0.74 |
| Agriculture, Forestry, Fishing and Hunting | 0.01 | 0.02 | 3.56 |
| Arts, Entertainment, and Recreation | 0.01 | 0.03 | 1.97 |
| Management of Companies and Enterprises | 0.01 | 0.02 | 2.47 |
| Other Services (except Public Administration) | 0.05 | 0.03 | 0.59 |
| Information | 0.01 | 0.02 | 1.24 |



Fig. 12: Percent of population in industry sectors.

shows the sensitivity of age groups to COVID-19 cases. Age groups 20-29, 65-79, and 80+ have larger Morris indices than other age groups, which indicates that the population from these groups is more sensitive to COVID-19 cases as changes in their population lead to greater changes in COVID-19 cases. The sensitivity of vaccination to COVID-19 cases is shown in Figure 11b. Vaccination is less distinguishable because all age subgroups converge at Morris indices at large delta values. The 40-49 and 0-19 age groups have slightly more impact on the cases than the 20-29 and 65-79 age groups in our experiments. This indicates that increasing the vaccination rate among these population groups will be equally impactful in

COVID-19 cases.

*D. Sensitivity of Vaccination by Industry Sectors*

We also studied the sensitivity of vaccination in our model by segmenting the population into the 18 majority represented industry sectors in Table X. To further rank the county-level population, Figure 12 represents the primary workplace setting across the 3142 US counties. The secondary Industry segments provide further details on common jobs in each of the 18 major segments [10], based on North American Industry Classification System (NAICS) from the Us Census Bureau. Note that the "Other Services" segment covers the rest of the industry jobs that cannot fit on the Table. Industry sectors
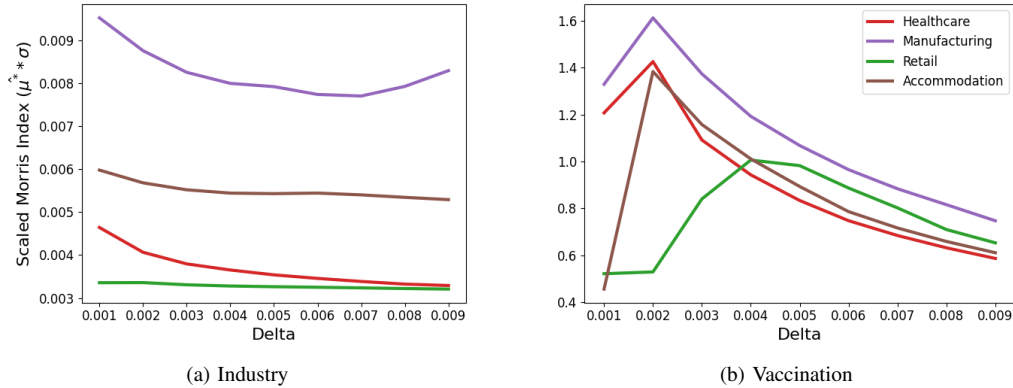
(a) Industry  (b) Vaccination

Fig. 13: Industry sectors and sensitivity results (the greater the value, the more important the feature sensitivity). We study features that measure a fraction of the county population by industry groups. We favor those with a large mean as they represent a larger proportion of the observed data. Quantitatively, we plot the scaled Morris index that multiplies the derivative with the standard deviation but only for those features with larger means (the top 4 industry sectors). Manufacture shows the most impact, followed by Accommodations in the Industry plot (left). The Vaccination plot (right) is less distinguishable between the sectors due to convergence at a large delta in the Morris indices.

in bold texts in Table X are studied in our experiments to indicate a higher population segment. The industry subgroups enable us to better understand the impact of COVID-19 in workplace settings. Looking at Fig. 13a, among the 4 industry sectors, increasing the Morris index in the Manufacturing sector seems to be the most sensitive to the infection, followed by Accommodation, Healthcare, and Retail.

## VII. RELATED WORK

### A. Statistical and Machine Learning Models

Many efforts have been made into COVID-19 forecasting using statistical learning, epidemiological, and deep learning models [18], [23], [24]. Different statistical and mathematical models such as Susceptible Infectious Recovery (SIR) and Susceptible Exposed Infectious Recover (SEIR), have been used to simulate the COVID-19 spread [24], [25]. These models provide useful insights, however, their performance is limited by the number of complex influencing factors and relationships they can capture, which is where the machine learning based models excel. Despite the difficulties in interpretability, deep learning models generally produce a better performance [23], [26] than the traditional machine learning and time series forecasting models, such as ARIMA [27], SVR [26], and XGBoost [28].

### B. Deep Learning Models and COVID-19 Forecasting

Deep learning has shown significant performance in time series forecasting [4] and has been widely used to predict COVID-19 infection. LSTM and Bi-LSTM based models often outperform other approaches in time-series forecasting [24], [26], [27] because RNN-based models are more suitable for time series data with Spatio-temporal sequences [29]. Furthermore, [18] found that Variational Auto Encoder (VAE) outperforms several other deep learning models in predicting

daily confirmed and recovered COVID-19 cases. AICov [30] and [16] used LSTM to predict COVID-19 in US counties. New research using the AI-based predictions also discusses 1) new virus variants are being discovered [31], 2) data quality and quantity used for model training are limited [25], [32], 3) ML-based models are not guaranteed to take socioeconomic, cultural and demographic factors into consideration while learning the data [32].

### C. Interpreting COVID-19 Forecasting Models

With the increasing use of deep learning models, their interpretation has gained increasing demand [33] to understand the model's decisions. DeepCOVID framework [34] used RNN with auto-regressive inputs to predict COVID-19 cases and interpret the contribution of input signals in prediction performance. DeepCOVIDNet [35] analyzed the features and their interactions in predicting the range of infected cases increases. Self-Adaptive Forecasting [2] is a novel method to adapt models on non-stationary time series data as well as giving interpretations with TFT.

## VIII. CONCLUSIONS AND FUTURE WORK

This paper uses deep learning to address the challenging problem of model feature sensitivity, which is critical in improving the interpretation of forecasting. With the COVID-19 data, we find that TFT performs well in learning temporal patterns from the data. The "self-attention" architecture sets up a strong global dependency structure. While it learns long-range dependencies, it is less effective in highly dynamic forecasting problems with non-stationary sequences, such as extreme or unknown events. However, the TFT model provides an excellent test environment enabling us to look into the COVID-19 feature sensitivity of high volatile time series and static variables. Combined with the Morris method, this sensitivity study enables us to look into stratifying the modeling

based on population subgroups such as industry and age and see how that affects COVID-19 responses to cases of infection and other important features at the community level.

The conceptualization of modeling individual features with the post hoc sensitivity analysis can apply to other time series models as the Morris method is efficient and generic. The advantages of using fewer features have enabled one to get insights into refined strategy optimizations. The Morris index and our novel individual multidimensional feature importance evaluation in this paper can contribute to autonomous feature engineering in forecasting and other similar challenges. COVID-19 infection prediction is a critical tool for policy-makers responding to a global health challenge with observed data changing in real-time. Future works can include analyzing the COVID-19 feature importance and descriptive accuracy in many health and financial problems.

### REFERENCES

[1] B. Wiegand, D. Klakow, and J. Vreeken, "Discovering interpretable data-to-sequence generators," in *Proceedings of the 36th AAAI Conference on Artificial Intelligence*, ser. AAAI '22, vol. 36, no. 4. AAAI Press, 2022, pp. 4237–4244.

[2] S. O. Arik, N. C. Yoder, and T. Pfister, "Self-adaptive forecasting for improved deep learning on non-stationary time-series," *arXiv preprint arXiv:2202.02403*, 2022.

[3] B. Peng, J. Li, S. Akkas, T. Araki, O. Yoshiyuki, and J. Qiu, "Rank position forecasting in car racing," in *2021 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. Los Alamitos, CA, USA: IEEE Computer Society, may 2021, pp. 724–733.

[4] B. Lim, S. Ö. Arık, N. Loeff, and T. Pfister, "Temporal fusion transformers for interpretable multi-horizon time series forecasting," *International Journal of Forecasting*, vol. 37, no. 4, pp. 1748–1764, 2021.

[5] USA Facts, "(us covid-19 cases and deaths)." [Online]. Available: https://usafacts.org/visualizations/coronavirus-covid-19-spread-map/

[6] Unacast, "Social Distancing Scoreboard." [Online]. Available: https://www.unacast.com/post/unacast-updates-social-distancing-scoreboard

[7] CDC, "COVID-19 Vaccinations in the United States Counties." [Online]. Available: https://data.cdc.gov/Vaccinations/COVID-19-Vaccinations-in-the-United-States-County/8xkx-amqh

[8] S. W. Marvel, J. S. House, M. Wheeler, K. Song, Y. Zhou, F. A. Wright, W. A. Chiu, I. Rusyn, A. Motsinger-Reif, and D. M. Reif, "The covid-19 pandemic vulnerability index (pvi) dashboard: Monitoring county-level vulnerability using visualization, statistical modeling, and machine learning," *medRxiv*, 2020.

[9] CDC, "Social Vulnerability Index," 2018. [Online]. Available: https://www.atsdr.cdc.gov/placeandhealth/svi/data_documentation_download.html

[10] "Secondary industry segments," (Date last accessed 23-August-2022). [Online]. Available: https://github.com/Data-ScienceHub/gpce-sensitivity

[11] A. Saltelli, M. Ratto, T. Andres, F. Campolongo, J. Cariboni, D. Gatelli, M. Saisana, and S. Tarantola, *Global sensitivity analysis: the primer*. John Wiley & Sons, 2008.

[12] G. Brunner, Y. Liu, D. Pascual, O. Richter, M. Ciaramita, and R. Wattenhofer, "On identifiability in transformers," *arXiv preprint arXiv:1908.04211*, 2019.

[13] M. D. Morris, "Factorial sampling plans for preliminary computational experiments," *Technometrics*, vol. 33, no. 2, pp. 161–174, 1991.

[14] F. Campolongo, J. Cariboni, and A. Saltelli, "An effective screening design for sensitivity analysis of large models," *Environmental modelling & software*, vol. 22, no. 10, pp. 1509–1518, 2007.

[15] PyTorch, "Temporal Fusion Transformer." [Online]. Available: https://pytorch-forecasting.readthedocs.io/en/stable/api/pytorch_forecasting.models.temporal_fusion_transformer.TemporalFusionTransformer.html

[16] M. D. Hssayeni, A. Chala, R. Dev, L. Xu, J. Shaw, B. Furht, and B. Ghoraani, "The forecast of covid-19 spread risk at the county level," *Journal of big data*, vol. 8, no. 1, pp. 1–16, 2021.

[17] J. Nossent and W. Bauwens, "Application of a normalized nash-sutcliffe efficiency to improve the accuracy of the sobol'sensitivity analysis of a hydrological model," in *EGU General Assembly Conference Abstracts*, 2012, p. 237.

[18] A. Zeroual, F. Harrou, A. Dairi, and Y. Sun, "Deep learning methods for forecasting covid-19 time-series data: A comparative study," *Chaos, Solitons & Fractals*, vol. 140, p. 110121, 2020.

[19] R. M. El-Shabasy, M. A. Nayel, M. M. Taher, R. Abdelmonem, K. R. Shoueir, and E. R. Kenawy, "Three waves changes, new variant strains, and vaccination effect against covid-19 pandemic," *International Journal of Biological Macromolecules*, vol. 204, pp. 161–168, 2022.

[20] US Census Bureau, "County population by characteristics: 2010-2020," https://www.census.gov/programs-surveys/popest/technical-documentation/research/evaluation-estimates/2020-evaluation-estimates/2010s-county-detail.html, 2020.

[21] J. M. McLaughlin, F. Khan, S. Pugh, D. L. Swerdlow, and L. Jodar, "County-level vaccination coverage and rates of covid-19 cases and deaths in the united states: An ecological analysis," *The Lancet Regional Health-Americas*, vol. 9, p. 100191, 2022.

[22] CBP, "County business patterns," 2020. [Online]. Available: https://www.census.gov/data/datasets/2020/econ/cbp/2020-cbp.html

[23] P. Wang, X. Zheng, G. Ai, D. Liu, and B. Zhu, "Time series prediction for the epidemic trends of covid-19 using the improved lstm deep learning method: Case studies in russia, peru and iran," *Chaos, Solitons & Fractals*, vol. 140, p. 110214, 2020.

[24] V. K. R. Chimmula and L. Zhang, "Time series forecasting of covid-19 transmission in canada using lstm networks," *Chaos, Solitons & Fractals*, vol. 135, p. 109864, 2020.

[25] J. S. Weitz, S. J. Beckett, A. R. Coenen, D. Demory, M. Dominguez-Mirazo, J. Dushoff, C.-Y. Leung, G. Li, A. Măgălie, S. W. Park *et al.*, "Modeling shield immunity to reduce covid-19 epidemic spread," *Nature medicine*, vol. 26, no. 6, pp. 849–854, 2020.

[26] F. Shahid, A. Zameer, and M. Muneeb, "Predictions for covid-19 with deep learning models of lstm, gru and bi-lstm," *Chaos, Solitons & Fractals*, vol. 140, p. 110212, 2020.

[27] İ. Kırbaş, A. Sözen, A. D. Tuncer, and F. Ş. Kazancıoğlu, "Comparative analysis and forecasting of covid-19 cases in various european countries with arima, narnn and lstm approaches," *Chaos, Solitons & Fractals*, vol. 138, p. 110015, 2020.

[28] J. Luo, Z. Zhang, Y. Fu, and F. Rao, "Time series prediction of covid-19 transmission in america using lstm and xgboost algorithms," *Results in Physics*, vol. 27, p. 104462, 2021.

[29] R. Chandra, A. Jain, and D. Singh Chauhan, "Deep learning via lstm models for covid-19 infection forecasting in india," *PloS one*, vol. 17, no. 1, p. e0262708, 2022.

[30] G. C. Fox, G. von Laszewski, F. Wang, and S. Pyne, "Aicov: An integrative deep learning framework for covid-19 forecasting with population covariates," *Journal of Data Science*, vol. 19, no. 2, pp. 293–313, 2021.

[31] E. A. Rashed and A. Hirata, "Infectivity upsurge by covid-19 viral variants in japan: Evidence from deep learning modeling," *International journal of environmental research and public health*, vol. 18, no. 15, p. 7799, 2021.

[32] C. Sáez, N. Romero, J. A. Conejero, and J. M. García-Gómez, "Potential limitations in covid-19 machine learning due to data source variability: A case study in the ncov2019 dataset," *Journal of the American Medical Informatics Association*, vol. 28, no. 2, pp. 360–364, 2021.

[33] B. Lim and S. Zohren, "Time-series forecasting with deep learning: a survey," *Philosophical Transactions of the Royal Society A*, vol. 379, no. 2194, p. 20200209, 2021.

[34] A. Rodriguez, A. Tabassum, J. Cui, J. Xie, J. Ho, P. Agarwal, B. Adhikari, and B. A. Prakash, "Deepcovid: An operational deep learning-driven framework for explainable real-time covid-19 forecasting," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 17, 2021, pp. 15 393–15 400.

[35] A. Ramchandani, C. Fan, and A. Mostafavi, "Deepcovidnet: An interpretable deep learning model for predictive surveillance of covid-19 using heterogeneous features and their interactions," *Ieee Access*, vol. 8, pp. 159 915–159 930, 2020.