Anomaly Detection in COVID-19 Time-Series Data

Hajar Homayouni $\,\cdot\,$ Indrakshi Ray $\,\cdot\,$ Sudipto Ghosh* $\,\cdot\,$ Shlok Gondalia $\,\cdot\,$ Michael G. Kahn

Abstract Anomaly detection and explanation in big volumes of real-world medical data, such as those pertaining to COVID-19, pose some challenges. First, we are dealing with time-series data. Typical time-series data describe behavior of a single object over time. In medical data, we are dealing with time-series data belonging to multiple entities. Thus, there may be multiple subsets of records such that records in each subset, which belong to a single entity are temporally dependent, but the records in different subsets are unrelated. Moreover, the records in a subset contain different types of attributes, some of which must be grouped in a particular manner to make the analysis meaningful. Anomaly detection techniques need to be customized for time-series data belonging to multiple entities. Second, anomaly detection techniques fail to explain the cause of outliers to the experts. This is critical for new diseases and pandemics where current knowledge is insufficient. We propose to address these issues by extending our existing work called IDEAL, which is an LSTM-Autoencoder based approach for data quality testing of sequential records, and provides explanations of constraint violations in a manner that is understandable to end-users. The extension (1) uses a novel two-level reshaping technique that splits COVID-19 datasets into multiple temporally-dependent sub-

H. Homayouni

Computer Science Dept., Colorado State University, Fort Collins CO 80523 E-mail: hajar.homayouni@colostate.edu

I. Ray

Computer Science Dept., Colorado State University, Fort Collins CO 80523 E-mail: iray@colostate.edu

S. Ghosh

Computer Science Dept., Colorado State University, Fort Collins CO 80523 E-mail: ghosh@colostate.edu

S. Gondalia

Computer Science Dept., Colorado State University, Fort Collins CO 80523 E-mail: shlok@rams.colostate.edu

M. Kahn

Anschutz Medical Campus, University of Colorado Denver, Aurora CO 80045 E-mail: michael.kahn@cuanschutz.edu

^{*} Corresponding author

sequences, and (2) adds a data visualization plot to further explain the anomalies and evaluate the level of abnormality of subsequences detected by IDEAL.

We performed two systematic evaluation studies for our anomalous subsequence detection. One study uses aggregate data, including the number of cases, deaths, recovered, and percentage of hospitalization rate, collected from a Covid tracking project, New York Times, and Johns Hopkins for the same time period. The other study uses COVID-19 patient medical records obtained from Anschutz Medical Center health data warehouse. The results are promising and indicate that our techniques can be used to detect anomalies in large volumes of real-world unlabeled data whose accuracy or validity is unknown.

Keywords Anomaly detection \cdot Data quality tests \cdot COVID-19 data \cdot Explainability \cdot LSTM-Autoencoder \cdot Time series

1 Introduction

Large amounts of data records are being collected from various sources over time to analyze the immediate and long-term impacts of COVID-19 on human health. Examples include the analysis of the impacts [46], diagnosis [25], treatments [17], and pre-symptom detection [38] of COVID-19 based on the available data collected from radiography, chest CT, chest X-ray, wearable devices, and COVID-19 tracking reports. Such data records are assumed to be accurate. However, the records may get corrupted in the non-trivial data collection and transformation processes. Anomalous data may lead to incorrect inferences and research findings. Thus, it is critical to automatically find inaccurate or anomalous data before doing any analysis and explain how the data is anomalous to the healthcare professionals.

Existing anomaly detection techniques for COVID-19 data (some based on machine learning) focus only on outbreak detection [28,26,52,11] in the COVID-19 tracking cases across the world. Machine learning approaches [49,30,48] have also been used for data quality assurance in other domains. Many of these techniques use supervised machine learning, which assumes the existence of labeled training data which in real-world is unavailable for new forms of outbreaks such as COVID-19. Moreover, outlier detection using machine learning fails to explain how records are anomalous to the domain experts. Finally, most works [15,14,16,44,21,23] identifying anomalous records in a dataset cannot be used on time-series data as anomalies may span multiple attributes and records in a sequence [35]. We aim to eliminate the above shortcomings by detecting anomalies in COVID-19 time-series data without having access to labeled data and explaining the anomalies to domain experts in a comprehensible manner.

This approach, called IDEAL, builds upon our previous work on data quality assessment approach [22] that uses an LSTM-Autoencoder [30] network to find anomalies in unsupervised data. Anomalies are data records or subsequences of data records whose behaviors (i.e., attribute values or change in the values over time) are significantly different from the majority of records and subsequences in a time-series dataset [28]. IDEAL automatically (1) discovers different types of constraints from the sequence data, (2) marks subsequences and records that violate the constraints as suspicious, and (3) explains the violations. IDEAL automatically generates three types of visualizations to explain the anomalies. The plot showing the suspiciousness score per attribute indicates which attributes make the subsequence anomalous. The second visualization uses decision trees to illustrate the violated constraints. The third plot compares a suspicious subsequence detected by IDEAL with normal subsequences belonging to the same dataset. The approach incorporates feedback from domain experts to improve the accuracy of constraint discovery and anomaly detection. We proposed an autocorrelation-based reshaping technique that automatically adjusts the LSTM-Autoencoder input window size based on how far the records are related to their past values. We evaluated the effectiveness of IDEAL using datasets from Yahoo servers [9], NASA Shuttle [4], and Colorado State University Energy Institute [2]. We demonstrated that IDEAL could detect previously known and injected anomalies in these datasets.

The above mentioned work needs to be extended for COVID-19 time-series data. The Yahoo servers [9] and NASA Shuttle [4] datasets that we previously used contain time-series data associated with a single entity; COVID-19 time-series data belongs to multiple entities (e.g., cases and deaths for states and counties, or lab test type and test name for different patients. Thus, in COVID-19 data there are multiple subsets of data each of which belongs to a single entity. Records in each subset are temporally dependent but they are unrelated to records in other subsets. Moreover, in each subset there are multiple grouping attributes (e.g., test type and test name) which requires the data to be preprocessed to make the results correct. We extend IDEAL by using a two-level reshaping approach to transform data into a shape that is suitable for analysis. This approach removes the restriction that all data records in a sequence dataset must be temporally dependent and are describing behaviors of the same object over time. Instead, IDEAL supports datasets in which a subset of records are temporally related to each other but are unrelated to the records from other subsets in the same dataset. For example, a health data store may contain medical records of multiple patients over time. The records of each patient are temporally dependent but independent from those of other patients.

One naïve solution may be to directly split the data based on grouping attribute(s) and generate multiple temporally-dependent subsequences. However, such a solution does not preserve associations among grouping attribute values. Consequently, IDEAL uses a pivoting-based approach to split data into multiple independent subsequences in a manner that preserves the associations among grouping attribute values.

In addition, IDEAL also offers an explanation of the anomalous behavior. It provides a data visualization plot that explains the level of abnormality of anomalous subsequences detected by the approach. This plot visualizes the data over time to help a domain expert understand the difference between the attribute values of a suspicious subsequence with those of other subsequences in the dataset. Such a plot draws attention to the anomalous subsequences; this is especially useful for large volumes of data where there is a lack of domain knowledge.

We conduct two types of studies to evaluate the anomaly detection effectiveness of IDEAL in the absence of domain knowledge. The first study validates the level of abnormality of an anomalous subsequence generated from a data source that is detected by IDEAL by comparing it with subsequences generated from other data sources that are reporting the same information. The data sources we use to conduct this study are COVID-19 tracking data collected from Johns Hopkins [5], New York Times [1], and COVID Tracking project [7] repositories. We demonstrate that IDEAL can detect anomalous subsequences which are indeed outliers when compared with other datasets reporting the same information.

The second study validates the level of abnormality of the suspicious subsequences by comparing the suspicious subsequences detected by IDEAL against other subsequences generated from the same dataset. Here the datasets correspond to a homogeneous population, i.e., a phenotype of people with the same values for personal (e.g., gender and age) and medical (e.g., diagnosis category, disease type, and medication type) attributes. We use COVID-19 medical data collected from the health data warehouse in Anschutz medical campus [3]. We demonstrate that IDEAL can detect abnormal subsequences from the datasets under test.

The contributions of this work are as follows:

- We propose a two-level reshaping technique to prepare data for training the LSTM-Autoencoder model. Thus, the preprocessing step that we develop allows IDEAL to be used on different types of time-series data, possibly grouped by different attributes, and belonging to multiple entities.
- We propose a data visualization plot to explain the level of abnormality of the subsequences detected by IDEAL. This helps the domain experts quickly identify the anomalous portions of data in large datasets.
- We propose systematic validation techniques based on a comparison between suspicious and other subsequences to demonstrate the anomaly detection effectiveness of IDEAL. Such a method is useful when there is a lack of labeled data or where there is insufficient domain knowledge.

The value of this work lies in automating the process of detecting and explaining potential anomalies that allow clinicians who have domain knowledge but lack data science skills to evaluate the effect of the level of abnormality and the seriousness of an anomaly on the clinical research question they are seeking to answer from the COVID-19 data. Due to the large number of investigators who intend to use the COVID-19 data, the use of the approach could potentially benefit a wide range of clinical investigators. This work can also be used for other domains that are analyzing large volumes of unlabeled time-series data that belong to multiple entities.

The rest of the paper is organized as follows. Section 2 describes the related work. Section 3 provides an overview of IDEAL. Section 4 describes how we handle time-series data belonging to multiple independent objects. Section 5 discusses anomaly interpretation in depth. Section 6 presents the evaluation of our approach. Section 7 concludes the paper and outlines directions for future work. Appendix A explains the architecture of the LSTM-Autoencoder network.

2 Related Work

The existing anomaly detection techniques in COVID-19 data focus only on outbreak detection [28,26,52,11] in the COVID-19 tracking cases across the world. Karadayi et al. [28] used a hybrid autoencoder network composed of a 3D convolutional neural network (CNN) and an autocorrelation based network for outbreak detection from spatio-temporal COVID-19 data provided by the Italian Department of Civil Protection. Jombart et al. [26] used linear regression, generalised

linear models (GLMs), and Bayesian regression to detect sudden changes in potential COVID-19 cases in England. However, there has been no focus on quality assurance of COVID-19 data used for various analysis.

Machine Learning-based techniques used for outlier detection in non-sequence data, such as Support Vector Machine (SVM) [15], Local Outlier Factor (LOF) [14], Isolation Forest (IF) [16], and Elliptic Envelope (EE) [44] have been used to detect anomalous records from time series data [42]. Such approaches do not consider temporal dependencies between data records and can only detect trivial out-of-range outliers.

Techniques that detect anomalous records from time-series data can be categorized as *decomposition* and *modeling* techniques. Decomposition techniques, suitable only for univariate time series, break a time series into level, trend, seasonality, and noise components and monitor the noise components to capture the anomalous records [24,33]. Modeling techniques represent a time series as a linear/non-linear function that associates each current value to its past values, predict the value of a record at a specific time, and report as anomalies those records whose prediction error falls outside a threshold. Stochastic modeling techniques, such as Moving Average (MA) [31], Autoregressive Integrated Moving Average (ARIMA) [37], and Holt-Winters (HW) [20] use statistical measures to calculate the correlation between the data records. These techniques assume that the time series is linear and follows a known statistical distribution, which make them inapplicable to many practical problems [10]. Machine learning modeling techniques support non-linear modeling, with no assumption about the distribution of the data [10]. Examples are Multi Layer Perceptrons (MLPs) [13], Long Short Term Memory (LSTM) [49], and Hierarchical Temporal Memory (HTM) [48]. Some of these techniques can model multivariate time-series. However, they produce complex equations, which are not human interpretable.

Existing techniques for anomalous sequence detection split the data into multiple subsequences, typically based on a fixed-size window [39] or an exhaustive brute-force approach [45]. Clustering-based anomalous sequence detection techniques extract subsequence features, such as trend and seasonality, and group the subsequences based on the similarities between their features. An anomalous subsequence is detected as the one that is distantly positioned within a cluster or is positioned in the smallest cluster. These approaches only detect anomalous sequences without determining the records and attributes that are the major causes of invalidity in each subsequence. Autoencoder-based techniques (1) take subsequences as input, (2) use an autoencoder network to reconstruct the subsequences, (3) assign invalidity scores based on the reconstruction errors to the subsequences, and (4) detect as anomalous those subsequences whose scores are greater than a threshold. These techniques can learn complex non-linear associations among the attributes in the time series but are not able to model the temporal dependencies among the records in the input subsequence. An LSTM-Autoencoder extends an autoencoder for time series data, and captures long-term temporal associations among data records in the form of complex equations that are not human interpretable. Our work aims to fill this gap by illustrating the cause of anomaly to the domain experts.

3 Approach

Figure 1 shows an overview of our approach. The input is in the form of data records and the output consists of a report showing subsequences of suspicious records accompanied with an explanation of the violated constraints. There are five components, namely, data preparation, constraint discovery, anomaly detection, anomaly interpretation, and anomaly inspection. These components form the basis of IDEAL [22] and are briefly described in the following paragraphs. Sections 4 and 5 describe how the data preparation and anomaly interpretation components are extended in this paper.



Fig. 1: IDEAL Overview [22]

Data Preparation This component prepares the data by transforming raw data into a form suitable for analysis. We used the one-hot encoding [50] method for preprocessing categorical attributes and the normalization [41] method for numeric attributes. Moreover, we proposed a systematic reshaping approach that uses autocorrelation [40] of the time-series attributes to enable the LSTM-Autoencoder network discover dependencies between highly correlated records. Note that, this step must be extended to handle COVID-19 data. The extensions are described in Section 4.

Constraint Discovery IDEAL uses an LSTM-Autoencoder, which is a sequenceto-sequence modeling technique [34] used to learn time series dependencies. An LSTM-Autoencoder can discover constraints involving long-term non-linear associations among multivariate time-series data records and attributes. The input and output to this network are fixed-size time series matrices. The network is composed of two hidden LSTM layers. The first LSTM layer functions as an encoder that investigates the dependencies from the input sequence and produces a complex hidden context. The second LSTM layer functions as a decoder that reconstructs the input sequence, based on the learned complex context and the previous output state. The difference between the original input and the reconstructed input is termed as the reconstruction error. Appendix A describes the architecture of the LSTM-Autoencoder network.

The LSTM-Autoencoder is an unsupervised technique that can potentially learn incorrect constraints from invalid data and generate false alarms. IDEAL uses an interactive learning approach that takes the expert's feedback through the anomaly inspection component to retrain the LSTM-Autoencoder model and improve its accuracy. Anomaly Detection This component detects suspicious subsequences and records that do not conform to the constraints represented by the trained model. Subsequence and records are assigned suspiciousness scores (s-scores), which are calculated based on the network reconstruction error and the record labels. The record label indicates the validity level of the record. If we start with an unlabeled dataset, the labels of all records are 0. The record label changes as we incorporate domain expert feedback in the subsequent iterations. Subsequences and records whose scores are greater than a threshold are flagged as suspicious. Using record labels in the definition of s-scores ensures that no valid subsequences or records are reported as suspicious in the retraining phase, thereby minimizing false alarms.

Anomaly Interpretation This component helps a domain expert interpret each suspicious subsequence by generating visualization plots of two types, namely, s-score per attribute and decision tree. The trained LSTM-Autoencoder model calculates the s-score per attribute. The higher the value of s-score, the more likely is the attribute to contribute to the invalidity of the subsequence. For each subsequence, IDEAL plots the s-score values for all the attributes in the subsequence. Moreover, IDEAL uses a decision tree [32] based technique called random forest [27] classifier to determine the constraints that are violated by each suspicious subsequence. For each attribute of the subsequence, a set of time series features, such as Mean, Max, and Curvature are extracted using Tsfeatures [43] CRAN library. Next, decision trees are generated using these features. The decision trees represent a set of if-then-else decision rules, which describe the constraints that identify sequences as valid or invalid based on their feature values. Note that, our anomaly interpretation component is extended in Section 5.

Anomaly Inspection This component takes domain expert feedback through a webbased user interface that uses check boxes for the expert to flag as faulty the subsequences that are actually anomalous. The feedback is incorporated to label the training data records as faulty or valid. The accuracy of constraint discovery is improved by adding the record label with four possible values (1: faulty, 0.5: suspicious, 0: unknown, and -1: valid) as a new attribute to the training dataset. This label is updated using domain expert feedback in every interaction. We redefine the reconstruction error of LSTM-Autoencoder based on the labels to minimize false alarms. The network is trained to minimize both the difference between the time series and its reconstruction, and the difference between the record labels in a time series and the labels predicted by the network.

4 Extension to Data Preparation

Analyzing COVID-19 time series data requires that the data be converted to a form suitable for analysis. In COVID-19 data, there are multiple subsets of data; records in each subset are temporally dependent but they are unrelated to records in other subsets. Such data must be prepared before being fed as input to sequential learning models, including the LSTM-based model used in this study, which assumes that all data records in an input sequence are temporally dependent and are describing behaviors of the same object over time.

A time series T is a sequence of d-dimensional records [29] described using the vector $T = \langle R_0, ..., R_{t-1} \rangle$, where $R_i = (a_i^0, ..., a_i^{d-1})$ is a record at time i, for $0 \le i \le t-1$ and a_i^j is the j^{th} attribute of the i^{th} record. A time series can be univariate (d=1) or multivariate (d>1) [19]. A univariate time series has one time-dependent attribute. For example, a univariate time series can consist of daily COVID-19 cases recorded sequentially over 24-hour increments. A multivariate time series is used to simultaneously capture the dynamic nature of multiple attributes. For example, a multivariate time series from a health data store can consist of multiple laboratory results of patients over time.

Reshaping is an essential data preparation step for sequential learning models [12,36]. This method reshapes the data to base the model computations at a time step t on a specified number of previous time steps. The number of previous time steps is known as window size.

Existing reshaping techniques use a single-level windowing approach, which assumes that all data records in a dataset are temporally dependent and are describing behaviors of a single object over time. For example, all the traffic data in the Yahoo Benchmark data store [9] are records related to a single server. The NASA Shuttle dataset [4] contains records of a single shuttle over time. However, real-world datasets including the ones used in this study typically contain records of multiple objects over time. For example, a COVID-19 tracking dataset can store case records of multiple states over time in the US. A medical dataset may contain clinical records for multiple patients over time. Each object (i.e., state in the COVID-19 dataset and *patient* in the medical dataset) has a unique id, which distinguishes the records concerning that object from the other records in the dataset. Single-level reshaping techniques cannot be used to split the data records into multiple subsequences in such datasets. These techniques may generate subsequences with temporally unrelated records, which can result in generating false alarms. For example, Table 1 shows a portion of records (i.e., for *Patient_ID*=1001 and *Patient_ID*=1005) in a medical dataset that stores patient weights over time.

$Patient_ID$	Timestamp	Weight
1001	6/1/2020	125.2
1001	7/1/2020	125.6
1005	7/1/2020	26.5
1001	8/1/2020	126.1
1005	8/1/2020	27

Table 1: Patient Weights Sequence

Splitting this dataset through single-level reshaping using window size equal to three with one record overlap results in two subsequences (Figure 2), which contain unrelated records.

However, the correct windowing must only contain temporally-related records of a single object (Figure 3). An anomaly detection technique may incorrectly detect subsequence 1 and 2 in Figure 2 as anomalous (i.e., two false positives) because of the sudden changes in the *Weight* values. Figure 3 shows the correct reshaping for our example.

Subsequence 1					
1001	6/1/2020	125.2			
1001	7/1/2020	125.6			
1005	7/1/2020	26.5			
:	Subsequence 2				
1005	7/1/2020	26.5			
1001	8/1/2020	126.1			
1005	8/1/2020	27			

Fig. 2: Incorrect Reshaping of Patient Records into Multiple Subsequences

Subsequence 1					
1001	6/1/2020	125.2			
1001	7/1/2020	125.6			
1001	8/1/2020	126.1			
Subsequence 2					
1005	7/1/2020	26.5			
1005	8/1/2020	27			

Fig. 3: Correct Reshaping of Patient Records into Multiple Subsequences

In this work, we propose a two-level reshaping technique to address the abovementioned issue. This technique (A) groups the time-series data based on domaindependent grouping attributes, and (B) splits the data records in each group using our systematic autocorrelation-based reshaping [22] approach.

4.1 Grouping Data Records

We split the data records into multiple temporally dependent groups. For this purpose, we (1) identify grouping attributes and their hierarchy, (2) concatenate the non-first-level grouping attributes into a new attribute, (3) pivot the new attribute into multiple temporal attributes, and (4) use the first-level grouping attribute to split the data records into multiple temporally-dependent subsequences.

(1) Identify grouping attributes. A grouping attribute is a categorical column by which we can group the dataset records into multiple temporally-dependent subsequences. A dataset may have one or more grouping attributes, which are domain-dependent. Figure 4a shows an example of a medical dataset of laboratory results for multiple patients over time. This dataset contains three levels of grouping attributes to describe the data records. The first-level grouping attribute (i.e., *Patient_ID*) indicates the objects in a dataset, each of which is represented by a unique Id. The second- to *h*-level grouping attributes indicate features about those objects. For example, in Figure 4a, the second-level grouping attribute (i.e., *Test_Type*) represents the type of laboratory test. Each patient can receive multiple types of test. The third-level grouping attribute (i.e., *Test_Name*) is the name of the laboratory test performed on the patients. Each *Test_Type* includes multiple *Test_Name*. We identify the domain-dependent grouping attributes with the

		1 st –level	2 nd –level	3 rd –level	
id	timestamp	Patient_ID	Test_Type	Test_Name	Value
1	19/9/2020	8	Blood	ALBUMIN	4.5
2	19/9/2020	8	Urinalysis	рН	4.7
3	19/9/2020	8	Urinalysis	Glucose	0.5
4	20/9/2020	8	Blood	ALBUMIN	9.0
5	21/9/2020	8	Blood	ALBUMIN	10.0
6	20/9/2020	9	Blood	ANIOGAR	3.9
7	21/9/2020	9	Blood	ANIOGAR	4.0
8	22/9/2020	9	Blood	ANIOGAR	11.1
9	22/9/2020	9	Blood	ALBUMIN	5.6
10	22/9/2020	9	Urinalysis	Protein	11
11	23/9/2020	9	Urinalysis	Protein	9

		1 st –level	2 nd –level	
id	timestamp	Patient_ID	2 nd –level Grouping Attribute	Value
1	19/9/2020	8	Blood-ALBUMIN	4.5
2	19/9/2020	8	Urinalysis-pH	4.7
3	19/9/2020	8	Urinalysis-Glucose	0.5
4	20/9/2020	8	Blood-ALBUMIN	9.0
5	21/9/2020	8	Blood-ALBUMIN	10.0
6	20/9/2020	9	Blood-ANIOGAR	3.9
7	21/9/2020	9	Blood-ANIOGAR	4.0
8	22/9/2020	9	Blood-ANIOGAR	11.1
9	22/9/2020	9	Blood-ALBUMIN	5.6
10	22/9/2020	9	Urinalysis-Protein	11
11	23/9/2020	9	Urinalysis-Protein	9

(a) Original Dataset with Multiple Grouping Attributes (b) Dataset after Concatenating Non-first-level Grouping

(b) Dataset after Concatenating Non-first-level Grouping Attributes

			1 st –level					
	id	timestamp	Patient_ID	Blood-ALBUMIN	Blood-ANIOGAR	Urinalysis-pH	Urinalysis-Glucose	Urinalysis-Protein
1	1	19/9/2020	8	4.5	Null	4.7	0.5	Null
G₁ —	2	20/9/2020	8	9.0	Null	Null	Null	Null
1	3	21/9/2020	8	10.0	Null	Null	Null	Null
	4	20/9/2020	9	Null	3.9	Null	Null	Null
~	5	21/9/2020	9	Null	4.0	Null	Null	Null
G ₂	6	22/9/2020	9	5.6	11.1	Null	Null	11
	7	23/9/2020	9	Null	Null	Null	Null	9

(c) Dataset after Pivoting Second-level Grouping Attribute

	id	timestamp	Patient_ID	Blood-ALBUMIN	Blood-ANIOGAR	Urinalysis-pH	Urinalysis-Glucose	Urinalysis-Protein	
	1	19/9/2020	8	4.5	Null	4.7	0.5	Null	Subsequence
G ₁ —	2	20/9/2020	8	9.0	Null	Null	Null	Null	Subsequence ₁
1	3	21/9/2020	8	10.0	Null	Null	Null	Null	
	4	20/9/2020	9	Null	3.9	Null	Null	Null	Subsequence
~	5	21/9/2020	9	Null	4.0	Null	Null	Null	Subsequence ₂
G ₂ –	6	22/9/2020	9	5.6	11.1	Null	Null	11	Subsoquonco
	7	23/9/2020	9	Null	Null	Null	Null	9	Subsequence ₃

(d) Dataset after Reshaping Using W = 2

Fig. 4: Splitting a Dataset with Multiple Grouping Attributes into Multiple Temporally-Dependent Subsequences

help from domain experts. In the future, we will use statistical autocorrelationbased [18] techniques to automatically identify the grouping attributes and their hierarchy from an input dataset.

(2) Concatenate non-first-level grouping attributes. Our approach converts all the non-first-level grouping attributes into a single dataset column to reduce

10

the complexity (i.e., dimensionality) of the problem. We call this new column the second-level grouping attribute. Figure 4b shows the new generated column from all the non-first-level grouping attributes.

At this step, we can use the first- and second-level attributes to group the data into multiple temporally dependent subsequences. However, using the resulting subsequences generated by this approach, time-series analysis techniques will not preserve the associations among the values of second-level attribute if any. For example, if there are associations among "Blood-Sugar" and "Blood-Pressure" of a patient, grouping at this stage would not preserve this association. To address this issue, our approach uses another step for pivoting the second-level attribute into multiple temporal attributes based on the attribute values.

(3) Pivot second-level grouping attribute. A pivoting query [6] converts all unique rows of an attribute into separate columns of their own, each of which contains a value specified as an input to the query. IDEAL pivots the second-level grouping attribute to generate multiple temporal attributes. The objective is to preserve the associations among grouping attribute values. Pivoting results in a smaller number of records in comparison to the original dataset.

Figure 4c shows how the pivoting process works in the example dataset. The second-level grouping attribute (i.e., a concatenation of *Test_Type* and *Test_Name*) is converted into five new attributes, each of which contains the corresponding test result stored in the *value* attribute. If a patient has not received a specific test at a specific time, the value of that test is set to Null.

(4) Group records by first-level grouping attribute. We use the first-level grouping attribute to categorize the data records in a sequence dataset into multiple groups G_i , $1 \le i \le m$, where m is the number of the distinct values of that attribute. Figure 4c shows the two groups of temporally-dependent records (i.e., G_1 and G_2) generated for the example dataset.

4.2 Autocorrelation-based Reshaping of Groups

For each group, IDEAL uses a systematic reshaping approach that we proposed in an earlier work [22] to split the data records in that group. This approach is based on the autocorrelation of the time series attributes to enable the LSTM-Autoencoder network discover dependencies between the records that are highly correlated. The input size is adjusted based on how far the records in a group are related to their past values. By feeding the LSTM-Autoencoder network with highly correlated records, this reshaping approach prevents the network from incorrectly discovering associations among non-correlated records. For each group G_i , IDEAL uses the autocorrelation-based approach to identify the window size w_i for that group. Autocorrelation is defined as the correlation of sequence data records with the records in the previous time steps, called lags [40]. An Autocorrelation Function (ACF [18]) at lag k for an attribute identifies to what extend the attribute is correlated to its k^{th} past value. IDEAL calculates ACF to identify the lags at which the attribute values are highly correlated to set the window size. As the LSTM-Autoencoder window size must be similar for all the data records in a dataset, IDEAL sets the final value of window size W to the smallest value of window sizes calculated for the groups (Equation 1). Finally, our approach reshapes the data records in groups based on the value of W.

$$W = Min(w_i), \ 1 \le i \le m \tag{1}$$

where m is the number of distinct values of the grouping attribute. Figure 4d shows how IDEAL splits the records in each group into multiple subsequences for W = 2.

5 Extension to Anomaly Interpretation

IDEAL uses an additional data visualization plot to explain the level of abnormality of suspicious subsequences. This plot visualizes attribute values for multiple groups (e.g., *state* in COVID-19 dataset and *patient* in medical dataset) over time. The visualization plot uses color-coded diagrams for every group to help a domain expert compare a suspicious subsequence with other subsequences from other groups in the dataset.



For each suspicious subsequence, IDEAL uses s-score per attribute values [22] to select the attribute with the highest suspiciousness score. Next, IDEAL plots values of that attribute for all the groups (i.e., G_i , $1 \leq i \leq m$) over time. The attribute values of the suspicious subsequence are represented by red points. Figures 5 and 6 shows the visualization plots generated for a suspicious subsequence detected from the laboratory results in the Anschutz medical data. Figure 5 shows that $e_{-}TOTGLOB$ attribute (i.e., total serum globulin) is the major cause of invalidity (i.e., attribute with the highest s-score value) in this subsequence. The data visualization plot in Figure 6 shows the values of the $e_{-}TOTGLOB$ attribute over time for this subsequence (in red) as well as other subsequences (in colors other than red) in the same dataset. We can visually observe from this figure how the values of $e_{-}TOTGLOB$ for this patient deviate from those of the majority of the patients in the dataset. As the value of this attribute is elevated in certain immunological diseases, this deviation can be caused by an immunological disease of the patient.

6 Evaluation

We evaluated the anomaly detection effectiveness of IDEAL using COVID-19 records from Johns Hopkins (JH) [5], New York Times (NT) [1], and Tracking project (T) [7] repositories. These publicly available datasets are updated daily and contain county- or state-level COVID-19 attributes. Wissel et al. [47] compared these datasets based on different factors, such as their data sources, collected attributes, region granularity, and frequency of updates. We used nine-month data from March 5th to November 11th, 2020 to evaluate one execution of IDEAL, which is an execution without the feedback loop. Moreover, we used four health datasets from the University of Colorado Anschutz medical campus [3] to evaluate the anomaly detection effectiveness of IDEAL. We used records of COVID-19-positive patients to evaluate one execution of IDEAL.

Current knowledge about the COVID-19 data attributes, pattern of spread, and distribution is insufficient as this is an unprecedented pandemic. We used two evaluation approaches to validate the suspicious subsequences detected from the COVID-19 data in the domain knowledge absence; these are (1) comparing suspicious subsequences detected by IDEAL from one data source to those from other independent sources that are recording values of the same records and attributes and (2) comparing the suspicious subsequences detected by IDEAL from a homogeneous population with the other subsequences in that population.

6.1 Comparing suspicious subsequences from different sources (Johns Hopkins (JH), New York Times (NT), and Tracking project (T))

The objective is to identify whether or not a suspicious subsequence is actually anomalous by comparing the suspicious subsequence from a data source with its equivalent subsequences from other independent COVID-19 data repositories. Two subsequences are equivalent if they contain records of same object (i.e., same grouping attribute value) and are observed during the same time period. For example, data records of the Alabama state collected from three sources of data during March to April 2020 form three equivalent subsequences. We formalized possible observations on equivalent subsequences to validate a suspicious subsequence based on whether (1) the same subsequence is detected by IDEAL as suspicious from all available sources of data or (2) the subsequence is detected as anomalous only in some of the available sources. We decided on whether or not each suspicious subsequence is actually anomalous based on a distance measure (i.e. mean square error (MSE)) between the attribute values of the suspicious subsequences collected from the other data sources.

For each suspicious subsequence s_i detected from the i^{th} source in a set of sources D (where |D| = n), we calculated the mean square error value between s_i and all its equivalent subsequences s_j from the other sources $(1 \le j \le n \text{ and } j \ne i)$. s_j can be either an undetected or a suspicious subsequence.

$$MSE(s_i, s_j) = \frac{1}{w} \sum_{k=1}^{w} (Normalized(A_i^f) - Normalized(A_j^f))^2$$
(2)

where A_i^f attribute is the major cause of invalidity in s_i , A_j^f is its equivalent attribute in s_j , and w is the window size. The following four cases describe how we validated a suspicious subsequence s_i based on the MSE value.

- (A) If attribute values of equivalent subsequences to s_i collected from all other sources of data are close to values of s_i , then all of those subsequences are either abnormal but valid, or anomalous detected from the same source. An abnormal subsequence can indicate signals of a COVID-19 outbreak [28].
 - If $\forall j \in \{1, ..., n\}, MSE(s_i, s_j) \leq Threshold$, then s_i and s_j are either – Abnormal but valid, or – Anomalous collected from the same source
- (B) If attribute values of equivalent subsequences to s_i collected from all other sources of data are far from values of s_i but close to one another, then s_i is anomalous:

If $\forall j,k \in \{1,...,n\}$ $(j,k \neq i), MSE(s_i,s_j) > Threshold$ and $MSE(s_j,s_k) \leq Threshold$, then s_i is: - Anomalous

(C) If attribute values of s_i are close to a subset of equivalent subsequences to s_i but far from another subset of equivalent subsequences, then all subsequences in the smaller subset are anomalous:

 $\begin{array}{l} \mbox{If } \forall j \in D_1, \ k \in D_2 \ (D_1 \cup D_2 = D, \ D_1 \cap D_2 = \varnothing, \ |D_1| < |D_2|), MSE(s_i,s_j) \leq \\ Threshold \ \mbox{and} \ MSE(s_i,s_k) > Threshold, \ \mbox{then:} \\ - \ s_i \ \mbox{and} \ s_j \ \mbox{are anomalous, and} \\ - \ s_k \ \mbox{is valid} \end{array}$

(D) If attribute values of s_i are close to a subset of equivalent subsequences to s_i but far from another subset of equivalent subsequences and the two subsets are of the same size, then all subsequences in both subsets are abnormal and need more investigations by a domain expert:

 $\begin{array}{ll} \text{If } \forall j \in D_1, \ k \in D_2 \ (D_1 \cup D_2 = D, \ D_1 \cap D_2 = \varnothing, \ |D_1| = |D_2|), MSE(s_i,s_j) \leq \\ Threshold \ \text{and} \ MSE(s_i,s_k) > Threshold, \ \text{then:} \end{array}$

 $-s_i, s_j$, and s_k are abnormal and need more investigation

Table 2: Datasets

Experime	ntDatasets	Attributes	
ID	(Sources)		
1	State-level data	Confirmed Cases (i.e.,	
from JH, NT, and		JH.Confirmed, NT.cases,	
	Т	T.positive) and Deaths	
2	State-level data	Recovered and Hospitaliza-	
	from JH and T	tion_Rate	

We ran two experiments. Table 2 shows the attributes from the COVID-19 datasets used for each experiment. Figure 7 shows data visualization plots for suspicious subsequences detected by IDEAL from datasets of the first experiment.

In this figure, each color represents data of a state. There are 50 plots for the 50 states of the US. The red plot represents the data of the suspicious subsequence. We used Figure 8 to validate the suspicious subsequences by comparing the attribute values of the suspicious subsequences with those of their equivalent subsequences from the first experiment datasets. These attributes are major causes of invalidity in each suspicious subsequence. In this experiment, the threshold T is set at 0.03 based on our observations of the values of MSE in these datasets.

Figure 7a. A suspicious subsequence s_{JH} was detected from JH in California data. An equivalent subsequence s_T was detected from T. The *Confirmed* attribute was the major cause of invalidity in these subsequences. The data visualization plot in Figure 7a shows how the *Confirmed* attribute values of the suspicious subsequence from California data in JH (red points) deviate from other subsequences from other states in the same source (i.e., JH). The constraint violations reported by the decision trees for this suspicious subsequence were over the *Minimum* and *Mean* features of the subsequence. In Figure 8a, $MSE(s_{JH}, s_T) < T$, $MSE(s_{JH}, s_{NT}) < T$, and $MSE(s_T, s_{NT}) < T$. This result indicates that in California, the numbers of confirmed cases over time reported by all three data sources for this suspicious subsequence are close to each other. Based on the case A, s_{JH} , s_{NT} , and s_T are either (i) abnormal but valid or (ii) anomalous data that have been obtained from the same source.

Figure 7b. A suspicious subsequence s_T was detected from T in New York data. Equivalent subsequences s_{JH} and s_{NT} were also detected from JH and NT. The *Deaths* attribute was the major cause of invalidity in these subsequences. The

16



(a) Values of *Confirmed* Attribute for US States over Time from JH. The Red Plot is a Suspicious Subsequence s_{JH} corresponding to California Data



(b) Values of death Attribute for US States over Time from T. The Red Plot is a Suspicious Subsequence s_T corresponding to New York Data



(c) Values of cases Attribute for US States over Time from NT. The Red Plot is a Suspicious Subsequence s_{NT} corresponding to Florida Data

Fig. 7: Visualization Plots for Suspicious Subsequences Detetced from Datasets of Experiment 1

data visualization plot in Figure 7b shows how the *Deaths* attribute values of the suspicious subsequence from New York data in T (red points) deviate from other subsequences from other states in the same source (i.e., T). The constraint violations reported by the decision trees for this suspicious subsequence were over the *Linearity* (i.e., strength of linearity, which is the sum of squared residuals of time-series from a linear autoregression) and *Burstiness* (i.e., ratio between the variance and the mean (Fano Factor) of time series) features of the subsequence.



Fig. 8: Actual Attribute Values in Suspicious Subsequences Detetced from Datasets of Experiment 1

In Figure 8b, $MSE(s_{JH}, s_T) > T$, $MSE(s_{NT}, s_T) > T$, and $MSE(s_{JH}, s_{NT}) < T$. This result indicates that the number of death cases in the New York state collected by the data source T was considerably less than that of the other two sources of data. Based on case B, s_T is anomalous and s_{JH}, s_{NT} are valid.



(a) Values of $Hospitalization_Rate$ Attribute for US States over Time from JH. The Red Plot is a Suspicious Subsequence s_{JH} corresponding to Kentucky Data



(b) Values of *Hospitalization_Rate* Attribute for US States over Time from T. The Red Plot is a Suspicious Subsequence s_T corresponding to Ohio Data



(c) Values of Hospitalization_Rate Attribute for US States over Time from T. The Red Plot is a Suspicious Subsequence s_T corresponding to Oregon Data

Fig. 9: Actual Attribute Values in Suspicious Subsequences Detetced from Datasets of Experiment 2

Figure 7c. A suspicious subsequence s_{NT} was detected only from NT in Florida data. The *Confirmed* attribute was the major cause of invalidity in this subsequence. The data visualization plot in Figure 7c shows how the *Confirmed* attribute values of the suspicious subsequence from Florida data in NT (red points) deviate from other subsequences from other states in the same source (i.e., NT).



Fig. 10: Actual Attribute Values in Suspicious Subsequences Detetced from Datasets of Experiment 2

The constraint violations reported by the decision trees for this suspicious subsequence were over the *Mean* and *Curvature* (i.e., strength of curvature, which is the amount by which a time series curve deviates from being a straight line and calculated based on the coefficients of an orthogonal quadratic regression) features of the subsequence. In this figure, $MSE(s_{JH}, s_T) < T$, $MSE(s_{JH}, s_{NT}) < T$, and $MSE(s_T, s_{NT}) < T$. This result indicates that in Florida, the numbers of confirmed cases over time reported by all three data sources for this suspicious subsequence are close to each other. Based on case A, s_{JH} , s_{NT} , s_T are either valid or anomalous collected from the same source.

Figure 9 shows data visualization plots for suspicious subsequences detected by IDEAL from datasets of the second experiment. We used Figure 10 to validate the suspicious subsequences by observing the actual values of attributes for the suspicious subsequences and their equivalent subsequences from the second experiment datasets. These attributes are major causes of invalidity in each suspicious subsequence. In this experiment threshold T is set at 0.0004 based on our observations on the values of MSE in these datasets.

Figure 9a. A suspicious subsequence s_{JH} was detected only from JH in Kentucky data. The *Hospitalization_Rate* attribute was the major cause of invalidity in this subsequence. The data visualization plot in Figure 9a shows how the *Hospitalization_Rate* attribute values of the suspicious subsequence (red points) from Kentucky data in JH deviate from other subsequences from other states in the same source (i.e., JH). The constraint violations reported by the decision trees for this suspicious subsequence were over the *Mean, Maximum*, and *Vchange* (i.e., maximum difference in variance between consecutive blocks in time series) features of the subsequence. In Figure 10a, $MSE(s_{JH}, s_T) > T$. This result indicates that the hospitalization rates reported by the two sources of data for the Kentucky state were considerably distinct. Based on case D, s_{JH} and s_T are abnormal subsequences that need more investigation.

Figure 9b. A suspicious subsequence s_T was detected only from T in Ohio data. The *Hospitalization_Rate* attribute was the major cause of invalidity in this subsequence. The data visualization plot in Figure 9b shows how the *Hospitalization_Rate* attribute values of the suspicious subsequence from Ohio data in T (red points) deviate from other subsequences from other states in the same source. The constraint violations reported by the decision trees for this suspicious subsequence were over the *Mean, Curvature* (i.e., strength of curvature, which is the amount by which a time series curve deviates from being a straight line and calculated based on the coefficients of an orthogonal quadratic regression), and *Highlowmu* (i.e., ratio between the means of data that is below and upper the global mean of time series) features of the subsequence. In Figure 10b, $MSE(s_{JH}, s_T) > T$. This result indicates that the hospitalization rates reported by the two sources of data for the Ohio state were considerably distinct. Based on case D, s_{JH} and s_T are abnormal subsequences that need more investigation.





Figure 9c. A suspicious subsequence s_T was detected only from T in Oregon data. The *Hospitalization_Rate* attribute was the major cause of invalidity in this subsequence. The data visualization plot in Figure 9c shows how the *Hospitalization_Rate* attribute values of the suspicious subsequence from Oregon data in T (red points) deviate from other subsequences from other states in the same source. The constraint violations reported by the decision trees for this suspicious subsequence were over the *Variance* and *Burstiness* (i.e., ratio between the variance and the mean (Fano Factor) of time series) features of the subsequence. In Figure 10c, $MSE(s_{JH}, s_T) > T$. This result indicates that the hospitalization rates reported by the two sources of data for the Oregon state were considerably distinct. Based on case D, s_{JH} and s_T are abnormal subsequences that need more investigation.

6.2 Comparing suspicious subsequences from a homogeneous population

The data of the patients of a homogeneous population should relatively look similar. We used this idea as a relative goal to evaluate the COVID-19 data in the absence of a domain expert. For this purpose, we extracted data of four homogeneous populations from Anschutz medical data store (Table 3). These datasets are results of joins of multiple tables (i.e., *Patient*, *Diagnosis*, and *Lab*) in the Anschutz health data warehouse. We fed each population data as an input dataset to the IDEAL tool to detect suspicious subsequences.

Table 3: Health Datasets

Datase	tDataset Name	#Records	#Attributes
ID			
1	COVID-positive with diabetes	770	103
2	COVID-positive females over 60	1174	103
3	COVID-positive with hypertension	1270	103
4	COVID-positive males over 60	1839	103

We validated the suspicious subsequences by visually observing the data visualization plots generated by IDEAL; we identified as actually abnormal (true positive) those suspicious subsequences of patients whose attribute values changing pattern over time are considerably different (i.e., visually observable) from other patients in their population. We identified as normal (false positive) those suspicious subsequences of patients whose attribute values changing pattern over time are not different (i.e., not visually observable) from other patients in their population.

Figure 11 shows a visualization plot generated by IDEAL for a suspicious subsequence detected by IDEAL from dataset ID=1. In this example, the $e_Meancorpusc3$ attribute is the major cause of suspiciousness of the subsequence. We can visually observe that the suspicious subsequence represented by red data points shows a considerable difference with other subsequences of the same population (i.e., COVID-positive with diabetes). As a result, this subsequence is a true positive. The constraint violations reported by the decision trees for this suspicious subsequence were over the *Mean* and *Variance* features of the subsequence. Figure 12 shows a visualization plot generated by IDEAL for a suspicious subsequence detected by IDEAL from dataset ID=3. In this example, the *Phart* attribute is the major cause of suspiciousness of the subsequence. We cannot visually observe a considerable difference between the suspicious subsequence represented by red data points and other subsequences of the same population (i.e., COVID-positive with hypertension). As a result, this subsequence is a false positive. The constraint violations reported by the decision trees for this suspicious subsequence were over the *Maximum*, *Curvature* (i.e., strength of curvature, which is the amount by which a time series curve deviates from being a straight line and calculated based on the coefficients of an orthogonal quadratic regression), and *Linearity* (i.e., strength of linearity, which is the sum of squared residuals of time series from a linear autoregression) features of the subsequence.

Table 4 shows number of true positives (TP), number of false positives (FP), $Precision = \frac{TP}{(TP+FP)}$ and total time (TT) it took to run the automated steps of IDEAL against each dataset under test. As the data is unlabeled, we cannot calculate the recall metric, which is based on the number of false negatives. The number of true positives and false positives are calculated based on our observation on the data visualization plots. It took between 58 to 108 seconds to run IDEAL against the datasets. IDEAL could detect between 1 to 4 abnormal subsequences in these datasets. The precision was between 75 to 100 percent.

Table 4: Results for Health Datasets

Dataset ID	TP	FP	Precision	TT (s)
1	4	1	0.80	58
2	3	1	0.75	76
3	2	0	1.00	64
4	1	0	1.00	108

7 Conclusions

We extended our previous data quality test approach to address the problem of anomaly detection in data pertaining to COVID-19. We (1) proposed a twolevel reshaping technique for data preparation, (2) added a data visualization plot for anomaly explanation, and (3) evaluated the approach against different COVID-19 datasets in the domain knowledge absence. We ran two experiments to validate the suspicious subsequences detected by IDEAL from Johns Hopkins, New York Times, and COVID-19 tracking project. We compared the attribute values of the suspicious subsequence detected from a source with those collected from other sources of data. IDEAL could find an anomalous subsequence in COVID-19 Tracking dataset in the number of deaths in the New York state. IDEAL could find three abnormal subsequences in the three sources, which need more investigation by domain experts.

We also evaluated the anomaly detection effectiveness of IDEAL using four health datasets from Anschutz medical campus. We compared a suspicious supsequence with other subsequences in a homogeneous population. IDEAL could detect ten abnormal subsequences in these datasets. In the future, we will evaluate the approach using other types of COVID-19 time series data. We plan to extend IDEAL to find anomalies in streaming COVID-19 data.

Acknowledgements This work was supported in part by funding from NSF under Award Numbers CNS 1650573, CNS 1822118, IIS 2027750, OAC 1931363, SecureNok, AFRL, Cable-Labs, Furuno Electric Company, NIST, and Google.

Conflict of Interest Statement

On behalf of all authors, the corresponding author states that there is no conflict of interest.

References

- 1. An Ongoing Repository of Data on Coronavirus Cases and Deaths in the U.S. https://github.com/nytimes/covid-19-data (Accessed 2020-08-25)
- 2. Energy Data. https://energy.colostate.edu/ (Accessed 2020-05-03)
- 3. HDC. http://www.ucdenver.edu/about/departments/healthdatacompass/ (Accessed 2020-11-20)
- 4. NASA Shuttle. https://archive.ics.uci.edu/ml/datasets/Statlog+(Shuttle) (Accessed 2020-05-15)
- 5. Novel Coronavirus COVID-19 (2019-nCoV) Data Repository by Johns Hopkins CSSE. https://github.com/CSSEGISandData/COVID-19 (Accessed 2020-08-25)
- Pivot Transformation. https://cloud.google.com/dataprep/docs/html/Pivot-Transform-57344645, (Accessed 2020-11-18)
- The COVID Tracking Project. https://github.com/COVID19Tracking (Accessed 2020-08-25)
- Understanding LSTM Networks, Recurrent Neural Networks. https://colah.github.io/ posts/2015-08-Understanding-LSTMs/ Accessed (29-03-2021)
- 9. Yahoo Server Traffic: A Benchmark Dataset for Time Series Anomaly Detection. https://yahooresearch.tumblr.com/post/114590420346/ a-benchmark-dataset-for-time-series-anomaly (Accessed 2020-07-01)
- Adhikari, R., Agrawal, R.K.: An Introductory Study on Time Series Modeling and Forecasting. LAP LAMBERT Academic Publishing (2013)
- Agbehadji, I.E., Awuzie, B.O., Ngowi, A.B., Millham, R.C.: Review of big data analytics, artificial intelligence and nature-inspired computing models towards accurate detection of covid-19 pandemic cases and contact tracing. International journal of environmental research and public health 17(15), 5330 (2020)
- Aljbali, S., Roy, K.: Anomaly Detection Using Bidirectional LSTM. In: K. Arai, S. Kapoor, R. Bhatia (eds.) Intelligent Systems and Applications, pp. 612–619. Springer International Publishing (2021)
- Bishop, C.M., Bishop, P.o.N.C.C.M.: Neural Networks for Pattern Recognition. Clarendon Press (1995)
- Breunig, M.M., Kriegel, H.P., Ng, R.T., Sander, J.: LOF: Identifying Density-Based Local Outliers. In: ACM SIGMOD International Conference on Management of Data, p. 93–104. Association for Computing Machinery (2000)
- Chen, Y., Wu, W.: Application of One-class Support Vector Machine to Quickly Identify Multivariate Anomalies from Geochemical Exploration Data. Geochemistry: Exploration, Environment, Analysis 17(3), 231–238 (2017)
- Cheng, Z., Zou, C., Dong, J.: Outlier Detection Using Isolation Forest and Local Outlier Factor. In: Conference on Research in Adaptive and Convergent Systems, p. 161–168. Association for Computing Machinery (2019)
- 17. Cortegiani, A., Ingoglia, G., Ippolito, M., Giarratano, A., Einav, S.: A systematic review on the efficacy and safety of chloroquine for the treatment of covid-19. Journal of critical care (2020)

- Finlay, R., Fung, T., Seneta, E.: Autocorrelation Functions. International statistical review 79(2), 255–271 (2011)
- Guo, T., Xu, Z., Yao, X., Chen, H., Aberer, K., Funaya, K.: Robust Online Time Series Prediction with Recurrent Neural Networks. In: IEEE International Conference on Data Science and Advanced Analytics, pp. 816–825 (2016)
- Hasani, Z., Jakimovski, B., Velinov, G., Kon-Popovska, M.: An Adaptive Anomaly Detection Algorithm for Periodic Data Streams", booktitle="Intelligent Data Engineering and Automated Learning. pp. 385–397. Springer International Publishing (2018)
- Homayouni, H., Ghosh, S., Ray, I.: ADQuaTe: An Automated Data Quality Test Approach for Constraint Discovery and Fault Detection. In: IEEE 20th International Conference on Information Reuse and Integration for Data Science, pp. 61–68. Los Angeles, CA (2019)
 Homayouni, H., Ghosh, S., Ray, I., Gondalia, S., Duggan, J., Kahn, M.G.: An
- Homayouni, H., Ghosh, S., Ray, I., Gondalia, S., Duggan, J., Kahn, M.G.: An autocorrelation-based lstm-autoencoder for anomaly detection on time-series data. In: 2020 IEEE International Conference on Big Data (Big Data), pp. 5068–5077 (2020)
- Homayouni, H., Ghosh, S., Ray, I., Kahn, M.: An Interactive Data Quality Test Approach for Constraint Discovery and Fault Detection. In: IEEE Big Data, pp. 200–205 (2019)
- Hyndman, R.J., Wang, E., Laptev, N.: Large-scale Unusual Time Series Detection. In: IEEE International Conference on Data Mining Workshop, pp. 1616–1619 (2015)
- Jain, R., Gupta, M., Taneja, S., Hemanth, D.J.: Deep learning based detection and analysis of covid-19 on chest x-ray images. Applied Intelligence pp. 1–11 (2020)
- Jombart, T., Ghozzi, S., Schumacher, D., Leclerc, Q., Jit, M., Flasche, S., Greaves, F., Ward, T., Eggo, R.M., Nightingale, E., et al.: Real-time monitoring of covid-19 dynamics using automated trend fitting and anomaly detection. medRxiv (2020)
- Kaminski, B., Jakubczyk, M., Szufel, P.: A Framework for Sensitivity Analysis of Decision Trees. Central European Journal of Operations Research 26(1), 135–159 (2018)
- Karadayi, Y., Aydin, M.N., Öğrencí, A.S.: Unsupervised anomaly detection in multivariate spatio-temporal data using deep learning: Early detection of covid-19 outbreak in italy. IEEE Access 8, 164155–164177 (2020). DOI 10.1109/ACCESS.2020.3022366
 Kieu, T., Yang, B., Jensen, C.S.: Outlier Detection for Multidimensional Time Series
- Kieu, T., Yang, B., Jensen, C.S.: Outlier Detection for Multidimensional Time Series Using Deep Neural Networks. In: 19th IEEE International Conference on Mobile Data Management, pp. 125–134 (2018)
- Kromkowski, P., Li, S., Zhao, W., Abraham, B., Osborne, A., Brown, D.E.: Evaluating Statistical Models for Network Traffic Anomaly Detection. In: Systems and Information Engineering Design Symposium, pp. 1–6 (2019)
- Kuster, C., Rezgui, Y., Mourshed, M.: Electrical Load Forecasting Models: A Critical Systematic Review. Sustainable Cities and Society 35, 257 – 270 (2017)
- 32. de Laat, P.B.: Algorithmic Decision-Making based on Machine Learning from Big Data: Can Transparency Restore Accountability. Philosophy & Technology 31(4), 525–541 (2018)
- Laptev, N., Amizadeh, S., Flint, I.: Generic and Scalable Framework for Automated Timeseries Anomaly Detection. In: 21th ACM International Conference on Knowledge Discovery and Data Mining, pp. 1939–1947 (2015)
- 34. Loganathan, G., Samarabandu, J., Wang, X.: Sequence to Sequence Pattern Learning Algorithm for Real-Time Anomaly Detection in Network Traffic. In: IEEE Canadian Conference on Electrical Computer Engineering, pp. 1–4 (2018)
- Lu, H., Liu, Y., Fei, Z., Guan, C.: An Outlier Detection Algorithm based on Cross-Correlation Analysis for Time Series Dataset. IEEE Access 6, 53593–53610 (2018)
- 36. Manaswi, N.K.: RNN and LSTM, pp. 115–126. Apress, Berkeley, CA (2018)
- Maçaira, P.M., Thomé, A.M.T., Oliveira, F.L.C., Ferrer, A.L.C.: Time Series Analysis with Explanatory Variables: A Systematic Literature Review. Environmental Modelling & Software 107, 199 – 209 (2018)
- Mishra, T., Wang, M., Metwally, A.A., Bogu, G.K., Brooks, A.W., Bahmani, A., Alavi, A., Celli, A., Higgs, E., Dagan-Rosenfeld, O., et al.: Pre-symptomatic detection of covid-19 from smartwatch data. Nature Biomedical Engineering pp. 1–13 (2020)
- Park, D., Hoshi, Y., Kemp, C.C.: A MultimodalAnomaly Detector for Robot-Assisted Feeding Using an LSTM-Based Variational Autoencoder. IEEE Robotics and Automation Letters 3(3), 1544–1551 (2018)
- 40. Park, K.I.: Fundamentals of Probability and Stochastic Processes with Applications to Communications, 1st edn. Springer Publishing Company, Incorporated (2017)
- Shalabi, L.A., Shaaban, Z.: Normalization as a Preprocessing Engine for Data Mining and the Approach of Preference Matrix. In: International Conference on Dependability of Computer Systems, pp. 207–214 (2006)

- Shriram, S., Sivasankar, E.: Anomaly Detection on Shuttle data using Unsupervised Learning Techniques. In: IEEE International Conference on Computational Intelligence and Knowledge Economy, pp. 221–225 (2019)
- Talagala, P.D., Hyndman, R.J., Smith-Miles, K., Kandanaarachchi, S., Munoz, M.A.: Anomaly Detection in Streaming Nonstationary Temporal Data. Journal of Computational and Graphical Statistics pp. 1–21 (2019)
- 44. Thomas, R., Judith, J.: Voting-Based Ensemble of Unsupervised Outlier Detectors. In: Advances in Communication Systems and Networks, pp. 501–511. Springer (2020)
- Wang, B., Wang, Z., Liu, L., Liu, D., Peng, X.: Data-driven Anomaly Detection for UAV Sensor Data Based on Deep Learning Prediction Model. In: 2019 Prognostics and System Health Management Conference, pp. 286–290 (2019)
- Wenham, C., Smith, J., Morgan, R.: Covid-19: the gendered impacts of the outbreak. The Lancet 395(10227), 846–848 (2020)
- 47. Wissel, B.D., Van Camp, P.J., Kouril, M., Weis, C., Glauser, T.A., White, P.S., Kohane, I.S., Dexheimer, J.W.: An interactive online dashboard for tracking COVID-19 in U.S. counties, cities, and states in real time. Journal of the American Medical Informatics Association 27(7), 1121–1125 (2020)
- Wu, J., Zeng, W., Yan, F.: Hierarchical Temporal Memory method for time-series-based anomaly detection. Neurocomputing 273, 535 – 546 (2018)
- Yu, Y., Si, X., Hu, C., Zhang, J.: A Review of Recurrent Neural Networks: LSTM Cells and Network Architectures. Neural Computation 31(7), 1235–1270 (2019)
- Zhang, W., Du, T., Wang, J.: Deep Learning over Multi-field Categorical Data. In: ECIR, pp. 45–57 (2016)
- Zhou, C., Paffenroth, R.C.: Anomaly Detection with Robust Deep Autoencoders. In: 23rd ACM International Conference on Knowledge Discovery and Data Mining, pp. 665–674 (2017)
- 52. Żhu, G., Li, J., Meng, Z., Yu, Y., Li, Y., Tang, X., Dong, Y., Sun, G., Zhou, R., Wang, H., et al.: Learning from large-scale wearable device data for predicting epidemics trend of covid-19. Discrete Dynamics in Nature and Society **2020** (2020)

Appendices

A LSTM-Autoencoder

A Long Short Term Network (LSTM) [49] is a Recurrent Neural Network (RNN) [8] that contains loops in its structure to allow information to persist and make network learn sequential dependencies among data records [49]. An RNN can be represented as multiple copies of a neural network, each passing a value to its successor. The original RNNs can only learn short-term dependencies among data records by using the recurrent feedback connections [19]. LSTMs extend RNNs by using specialized gates and memory cells in their neuron structure to learn long-term dependencies. The computational units (neurons) of an LSTM are called *memory cells*. An LSTM has the ability to remove or add information to the memory cell state by using *gates*. The gates are defined as weighted functions that govern information flow in the memory cells. The gates are composed of a *sigmoid layer* and a *point-wise operation* to optionally let information through. The sigmoid layer outputs a number between zero (to let nothing through) and one (to let everything through). There are three types of gates, namely, *forget, input*, and *output*.

 Forget gate: Decides what information to discard from the memory cell. Equation 3 shows the mathematical representation of the forget gate.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \tag{3}$$

where W_f is the connection weight between the inputs $(h_{t-1} \text{ and } x_t)$ and the sigmoid layer; b_f is the bias term and σ is the sigmoid activation function. In this gate, $f_t = 1$ means that completely keep the information and $f_t = 0$ means that completely get rid of the information.

- *Input gate*: Decides which values to be used from the network input to update the memory state. Equation 4 shows the mathematical representation of the input gate.

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \tag{4}$$

where C_t is the new memory cell state and C_{t-1} is the old cell state, which is multiplied by f_t to forget the information decided by the forget gate; \tilde{C}_t is the new candidate value for the memory state, which is scaled by i_t as how much the gate decides to update the state value.

- Output gate: Decides what to output based on the input and the memory state. Equation 5 shows the mathematical representation of the output gate. This gate pushes the cell state values between -1 and 1 by using a hyperbolic tangent function and multiplies it by the output of its sigmoid layer to decide which parts of the input and the cell state to output.

$$o_t = \sigma(W_o.[h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * tanh(C_t)$$
(5)

An autoencoder is an unsupervised deep neural network that discovers constraints in the unlabeled input data. An autoencoder is composed of an *encoder* and a *decoder*. The encoder compresses the data from the input layer into a short representation, which is a non-linear combination of the input elements. The decoder decompresses this representation into a new representation that closely matches the original data. The network is trained to minimize the reconstruction error (RE), which is the average squared distance between the original data and its reconstruction [51].

An LSTM-Autoencoder [29] is an extension of an autoencoder for time-series data using an encoder-decoder LSTM architecture. An LSTM-Autoencoder can capture the temporal dependencies among the input records by using LSTM networks as the layers of the autoencoder network.



Fig. 13: An LSTM-Autoencoder Network

Figure 13 shows the LSTM-Autoencoder architecture. The input and output are fixedsize time series matrices. $X_{i,j} = [x_{i,j}^0, ..., x_{i,j}^{d-1}]$ is the j^{th} record with d attributes, T_i is the i^{th} time series that contains w records, and w is the window size. The network output has the same dimensionality as the network input. The network is composed of two hidden layers that are LSTMs with d' units. The first LSTM layer functions as an encoder that investigates the dependencies from the input sequence and produces a complex hidden context (i.e., d' encoded time series features, where the value of d' depends on the underlying encoding used by the autoencoder). The second LSTM layer functions as a decoder that produces the output sequence, based on the learned complex context and the previous output state. The TimeDistributed layer is used to process the output from the LSTM hidden layer. This layer is a dense (fully-connected) wrapper layer that makes the network return a sequence with shape (d * w). The reconstruction error for this network is defined as follows [51]:

$$RE = \frac{1}{m} \sum_{i=1}^{m} (T'_i - T_i)^2 \tag{6}$$

where T_i and T_i^\prime are the i^{th} network input and output and m is the total number of subsequences.