

Anomaly Detection for High-Dimensional Data Using Large Deviations Principle

Sreelekha Guggilam
sreelekh@buffalo.edu
University at Buffalo
Buffalo, NY, USA

Varun Chandola
chandola@buffalo.edu
University at Buffalo
Buffalo, NY, USA

Abani Patra
abani.patra@tufts.edu
Tufts University
Boston, MA, USA

ABSTRACT

Most current anomaly detection methods suffer from the *curse of dimensionality* when dealing with high-dimensional data. We propose an anomaly detection algorithm that can scale to high-dimensional data using concepts from the theory of *large deviations*. The proposed *Large Deviations Anomaly Detection* (LAD) algorithm is shown to outperform state of art anomaly detection methods on a variety of large and high-dimensional benchmark data sets. Exploiting the ability of the algorithm to scale to high-dimensional data, we propose an online anomaly detection method to identify anomalies in a collection of multivariate time series. We demonstrate the applicability of the online algorithm in identifying counties in the United States with anomalous trends in terms of COVID-19 related cases and deaths. Several of the identified anomalous counties correlate with counties with documented poor response to the COVID pandemic.

CCS CONCEPTS

• **Computing methodologies** → **Anomaly detection**.

KEYWORDS

Large deviations, anomaly detection, high-dimensional data, multivariate time series

1 INTRODUCTION

Anomaly detection has been extensively studied over many decades across many domains [9, 18]. Among the most useful applications of anomaly detection is to simultaneously monitor multiple systems' behaviors and identify the system that exhibits anomalous behavior due to external or internal stress factors. For instance, consider the example of the COVID-19 infection data. Studying the confirmed case and death trends across various countries, states or counties could highlight and identify the most (or least) significant public policies. One possible approach to study the data could be to monitor each time series [8, 20, 30] and identify sudden outbreaks or significant causal events. However, such methods study each time series individually and cannot not be used to detect the gradual divergence from the normal trends or initial signs of such drift.

An alternate approach is to analyze each time series in the context of a collection of time series, which can reveal anomalies beyond sudden and significant events, such as anomalous trends and gradual drifts. Such methods typically require an appropriate similarity measure [16]. Through appropriate combination with state-of-the-art similarity-based models, these methods can identify potential anomalous time series and cluster similar trends. Implementing such methods in a time varying setting could even help

detect change points or anomalous events in individual time series as well as identifying anomalous time series [5, 31]. However, these methods are typically unable to scale to long time series [4, 31].

In this paper, we propose a new anomaly detection algorithm called *Large deviations Anomaly Detection* (LAD), for large/high-dimensional data and multivariate time series data. LAD uses the rate function from *large deviations principle* (LDP) [14, 27, 28] to deduce anomaly scores for the underlying data. Core ideas for the algorithm are inspired from large deviation theory's projection theorem that allow better handling of high dimensional data. Unlike most high dimensional anomaly detection models, LAD does not incorporate feature selection or dimensionality reduction, which makes it ideal to study multiple time series in an online mode. The intuition behind the LAD model allows it to naturally segregate the anomalous observations at each time step while comparing multiple multivariate time series simultaneously. The key contributions of this paper are following:

- (1) We propose the *Large deviations Anomaly Detection* (LAD) algorithm, a novel and highly scalable LDP based methodology, for scoring based anomaly detection.
- (2) The proposed LAD model is capable of analyzing large and high dimensional datasets without additional dimensionality reduction procedures thereby allowing more accurate and cost effective anomaly detection.
- (3) An online extension of the LAD model is presented to detect anomalies in an multivariate time series database using an evolving anomaly score for each time series. The anomaly score varies with time and can be used to track developing anomalous behavior.
- (4) We perform an empirical study on publicly available anomaly detection benchmark datasets to analyze robustness and performance of the proposed method on high dimensional and large datasets.
- (5) We present a detailed analysis of COVID-19 trends for US counties where we identify counties with anomalous behavior (See Figure 1 for an illustration).

The rest of this document is organized as follows. Section 2 provides an overview of relevant existing methods for anomaly detection. Section 3 is a short background on underlying large deviations theory motivating LAD. Section 4 details our LAD model for detecting unsupervised anomalies in multivariate time series. Section 5 describes the experiments and demonstrate the state-of-the-art performance of our method. Section 6 concludes the paper and sketches direction for possible future work.

¹In early November, these counties in North Dakota were exhibiting infection rates that were six times the national rate - <https://www.washingtonpost.com/opinions/2020/11/06/north-dakota-covid-19-cases/>

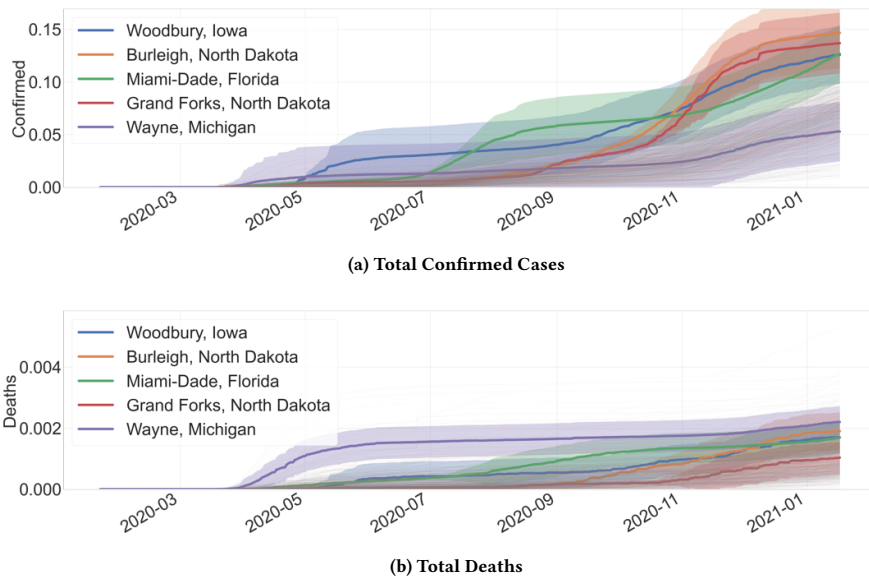


Figure 1: Top 5 anomalous counties identified by the proposed LAD algorithm based on the daily multivariate time-series, consisting of cumulative COVID-19 per-capita infections and deaths. At any time-instance, the algorithm analyzes the bi-variate time series for all the counties to identify anomalies. The time-series for the non-anomalous counties are plotted (light-gray) in the background for reference. For the counties in North Dakota (Burleigh and Grand Forks), the number of confirmed cases (*top*), and the sharp rise in November 2020, is the primary cause for anomaly¹. On the other hand, Wayne County in Michigan was identified as anomalous primarily because of its abnormally high death rate, especially when compared to the relatively moderate confirmed infection rate.

2 RELATED WORK

In this section, we provide a brief overview of relevant anomaly detection methods which have been proposed for high-dimensional data and for multivariate time-series data. We also discuss other works that have used the large deviations principle for detecting anomalies.

A large body of research exists on studying anomalies in high dimensional data [1, 3] but challenges remain. Many anomaly detection algorithms use dimensionality reduction techniques as a pre-processing step to anomaly detection. However, many high dimensional anomalies can only be detected in high dimensional problem settings and dimensionality reduction in such settings can lead to false negatives. Many methods exist that identify anomalies on high-dimensional data without dimensional reduction or feature selection, e.g. by using distance metrics. *Elliptic Envelope* (EE) [25] fits an ellipse around data centers by fitting a robust covariance estimates. *Isolation Forest* (I-Forest) [19] uses recursive partitioning by random feature selection and isolating outlier observations. *k nearest neighbor outlier detection* (kNN) [23] uses distance from nearest neighbor to get anomaly scores. *local outlier factor* (LOF) [7] uses deviation in local densities with respect to its neighbors to detect anomalies. *k-means*-- [12] method uses distance from nearest cluster centers to jointly perform clustering and anomaly detection. *Concentration Free Outlier Factor* (CFOF) [2] uses a “reverse nearest neighbor-based score” which measures the number of nearest neighbors required for a point to have a set proportion of data within its envelope. In particular, methods like I-Forest and CFOF are targeted towards anomaly detection in high dimensional datasets.

In most settings, real time detection of anomalies is needed to dispatch necessary preventive measures for damage control. Such problem formulation requires collectively monitoring a high dimensional time series database to identify anomalies in real time. Recently, large deviations theory has been widely applied in the fields of climate models [13], statistical mechanics [26], networks [22], etc. Specially for analysis of time series, the theory of large deviations has proven to be of great interest over recent decades [6, 21]. However, these methods are data specific, often study individual time series and are difficult to generalize to other areas of research.

Anomaly detection for time series have been extensively explored in the literature [17], though most focus has been on identifying anomalous events in a single time-series. While, the task of detecting anomalous time series in a collection of time series has been studied in the past [10, 11, 29], most of these works have focused on univariate time series and have not shown to scale to long time series data. Our proposed method addresses this issue by using the large deviation principle.

3 LARGE DEVIATION PRINCIPLE

Large deviations theory provides techniques to derive the probability of rare events² that have an asymptotically exact exponential approximation [14, 27, 28]. In this section, we briefly go over the large deviation theory and different ways to generate the rate functions required for the large deviations principle.

²In our context, these rare events include outlier/anomalous behaviors.

The key concept of this theory is the Large Deviations Principle (LDP). The principle describes the exponential decay of the probabilities for the mean of random variables. The rate of decay is characterized by the rate function \mathcal{I} . The theorem is detailed below:

THEOREM 3.1. *A family of probability measures $\{\mu_\epsilon\}_{\epsilon>0}$ on a Polish space \mathcal{X} is said to satisfy large deviation principle (LDP) with the rate function $\mathcal{I} : \mathcal{X} \rightarrow [0, \infty]$ if:*

- (1) \mathcal{I} has compact level sets and is not identically infinite
- (2) $\liminf_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(O) \geq -\mathcal{I}(O) \quad \forall O \subseteq \mathcal{X}$ open sets
- (3) $\limsup_{\epsilon \rightarrow 0} \epsilon \log \mu_\epsilon(C) \leq -\mathcal{I}(C) \quad \forall C \subseteq \mathcal{X}$ closed sets

where, $\mathcal{I}(S) = \inf_{x \in S} \mathcal{I}(x)$, $S \subseteq \mathcal{X}$

To implement LDP on known data with known distributions, it is important to decipher the rate function \mathcal{I} . Cramer's Theorem provides the relation between the rate function \mathcal{I} and the logarithmic moment generating function Λ .

Definition 3.2. The logarithmic moment generating function of a random variable X is defined as

$$\Lambda(t) = \log E[\exp(tX)] \quad (1)$$

THEOREM 3.3 (CRAMER'S THEOREM). *Let X_1, X_2, \dots, X_n be a sequence of iid real random variables with finite logarithmic moment generating function, e.g. $\Lambda(t) < \infty$ for all $t \in \mathbb{R}$. Then the law for the empirical average satisfies the large deviations principle with rate $\epsilon = 1/n$ and rate function given by*

$$\mathcal{I}(x) := \sup_{t \in \mathbb{R}} (tx - \Lambda(t)) \quad \forall t \in \mathbb{R} \quad (2)$$

Thus, we get,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \left(P \left(\sum_{i=1}^n X_i \geq nx \right) \right) = -\mathcal{I}(x), \quad \forall x > E[X_1] \quad (3)$$

For more complex distributions, identifying the rate function using logarithmic moment generating function can be challenging. Many methods like contraction principle and exponential tilting exist that extend rate functions from one topological space that satisfies LDP to the topological spaces of interest[14]. For our work, we are interested in the Dawson-Gärtner Projective LDP, that generates the rate function using nested family of projections.

THEOREM 3.4. Dawson-Gärtner Projective LDP: *Let $\{\pi^N\}_{N \in \mathbb{N}}$ be a nested family of projections acting on \mathcal{X} s.t. $\cup_{N \in \mathbb{N}} \pi^N$ is the identity. Let $\mathcal{X}^N = \pi^N \mathcal{X}$ and $\mu_\epsilon^N = \mu_0 \circ (\pi^N)^{-1}$, $N \in \mathbb{N}$. If $\forall N \in \mathbb{N}$, the family $\{\mu_\epsilon^N\}_{\epsilon>0}$ satisfies the LDP on \mathcal{X}^N with rate function \mathcal{I}^N , then $\{\mu_\epsilon\}_{\epsilon>0}$ satisfies the LDP with rate function \mathcal{I} given by,*

$$\mathcal{I}(x) = \sup_{N \in \mathbb{N}} \mathcal{I}^N(\pi^N x) \quad x \in \mathcal{X}$$

Since $\mathcal{I}^N(y) = \inf_{\{x \in \mathcal{X} | \pi^N(x)=y\}} \mathcal{I}(x)$, $y \in \mathcal{Y}$, the supremum defining \mathcal{I} is monotone in N because projections are nested.

The theorem allows extending the rate function from a lower projection to higher projection space. The implementation of this theorem in LAD model is discussed in Section 4.

4 METHODOLOGY

Consider the case of multivariate time series data. Let $\{t_n\}_{n=1}^N$ be a set of multivariate time series datasets where $t_n = (t_{n,1}, \dots, t_{n,T})$ is a time series of length T and each $t_{n,t}$ has d attributes. The motivation is to identify anomalous t_n that diverge significantly from the non-anomalous counter parts at a given or multiple time steps.

The main challenge is to design a score for individual time series that evolves in a temporal setting as well as enables tracking the initial time of deviation as well as the scale of deviation from the normal trend.

As shown in following sections, our model addresses the problem through the use of rate functions derived from large deviations principle. We use the Dawson-Gärtner Projective LDP (See Section 4.2) for projecting the rate function function to a low dimensional setting while preserving anomalous instances.

The extension to temporal data (See Section 4.3) is done by collectively studying each time series data as one observation.

4.1 Large Deviations for Anomaly Detection

Our approach uses a direct implementation of LDP to derive the rate function values for each observation. As the theory focuses on extremely rare events, the raw probabilities associated with them are usually very small [14, 27, 28]. However, the LDP provides a rate function that is useful as a scoring metric for our LAD model.

Consider a dataset X of size n . Let $a = \{a_1, \dots, a_n\}$ and $I = \{I_1, \dots, I_n\}$ be anomaly score and anomaly label vectors for the observations respectively such that $a_i \in [0, 1]$ and $I_i \in \{0, 1\}$ $\forall i \in \{1, 2, \dots, n\}$.

By large deviations principle, we know that for a given dataset X of size n , $P(\bar{X} = p) \approx e^{-n\mathcal{I}(p)}$. Assuming that the underlying data is standard Gaussian distribution with mean 0 and variance 1, we can use the rate function for Gaussian data where $\mathcal{I}(p) = \frac{p^2}{2}$. Then the resulting probability that the sample mean is p is given by:

$$P(\bar{X} = p) \approx e^{-n\frac{p^2}{2}} \quad (4)$$

Now, in presence of an anomalous observation x_a , the sample mean is shifted by approximately x_a/n for large n . Thus, the probability of the shifted mean being the true mean is given by,

$$P(\bar{X} = x_a/n) \approx e^{-\frac{x_a^2}{2n}} \quad (5)$$

However, for large n and $|x_a| \ll 1$, the above probabilities decay exponentially which significantly reduces their effectiveness for anomaly detection. Thus, we use $\frac{x_a^2}{2n}$ as anomaly score for our model. Thus generalizing this, the anomaly score for each individual observation is given by:

$$a_i = n\mathcal{I}(x_i) \quad \forall i \in \{1, 2, \dots, n\} \quad (6)$$

4.2 LDP for High Dimensional Data

High dimensional data pose significant challenges to anomaly detection. Presence of redundant or irrelevant features act as noise making anomaly detection difficult. However, dimensionality reduction can impact anomalies that arise from less significant features

of the datasets. To address this, we use the Dawson-Gärtner Projective theorem in LAD model to compute the rate function for high dimensional data. The theorem records the maximum value across all projections which preserves the anomaly score making it optimal to detect anomalies in high dimensional data. The model algorithm is presented in Algorithm 1.

Algorithm 1: Algorithm 1: LAD Model

Input: Dataset X of size (n, d) , number of iterations N_{iter} , threshold th .

Output: Anomaly score a

Initialization: Set initial anomaly score and labels a and I to zero vectors and, entropy matrix $E = 0_{(n,d)}$ where $0_{(n,d)}$ is a zero matrix of size (n, d) .

for each $s \rightarrow 1$ **to** N_{iter} **do**

- (1) Subset $X_{sub} = X[I_i == 0]$
 - (2) $X_{normalized}[:, d_i] = \frac{X[:, d_i] - X_{sub}[:, d_i]}{cov(X_{sub}[:, d_i])}, \quad \forall d_i \in \{1, \dots, d\}$
 - (3) $E[i, :] = -X_{normalized}[i]^2/2n, \quad \forall i$
 - (4) $a_i = -\max(E[i, :])$
 - (5) $a = \frac{a - \min(a)}{\max(a) - \min(a)}$
 - (6) $th = \min(th, \text{quantile}(a, 0.95))$
 - (7) $I_i = 1$ if $a_i > th, \quad \forall i$
-

4.3 LAD for Time Series Data

The definition of an anomaly is often contingent on the data and the problem statement. Broadly, time series anomalies can be categorized to two groups [10]:

- (1) **Divergent trends/Process anomalies:** Time series with divergent trends that last for significant time periods fall into this group. Here, one can argue that generative process of such time series could be different from the rest of the non-anomalous counterparts.
- (2) **Subsequence anomalies:** Such time series have temporally sudden fluctuations or deviations from expected behavior which can be deemed as anomalous. These anomalies occur as a subsequence of sudden spikes or fatigues in a time series of relatively non-anomalous trend.

The online extension of the LAD model is designed to capture anomalous behavior at each time step. Based on the mode of analysis of the temporal anomaly scores, one can identify both divergent trends and subsequence anomalies. In this paper, we focus on the divergent trends (or process anomalies). In particular, we try to look at the anomalous trends in COVID-19 cases and deaths in US counties. Studies to collectively identify divergent trends and subsequence anomalies is being considered as a prospective future work.

In this section, we present an extension of the LAD model to multivariate time series data. Here, we wish to preserve the temporal dependency as well as dependency across different features of the time series. Thus, as shown in Algorithm 2, a horizontal stacking of the data is performed. This allows collective study of temporal and non-temporal features. To preserve temporal dependency, the

anomaly scores and labels are carried on to next time step where the labels are then re-evaluated.

Algorithm 2: Algorithm 2: LAD for Time series anomaly detection

Input: Time series dataset $\{t_n\}_{n=1}^N$ of size (N, T, d) , number of iterations N_{iter} , threshold th , window w .

Output: An array of temporal anomaly scores a , an array of temporal anomaly labels I

Initialization: Set initial anomaly score and labels a and I to zero matrices of size (N, T) and, entropy matrix E to a zero matrix of size (N, T, d) .

for each $t \rightarrow 1$ **to** T **do**

$X = \text{hstack}(t_{n,t}^-)$ where $t_{n,t}^- = \{t_{n,t-w}, \dots, t_{n,t}\}$

$I[i, t] = I[i, t-1]$

$a[:, t] = a[:, t-1]$

for each $s \rightarrow 1$ **to** N_{iter} **do**

(1) Subset non-anomalous time series

$X_{sub} = \{X[i, :] | I[i, t] == 0, \forall i\}$

(2) $X_{normalized}[:, d_i] = \frac{X[:, d_i] - X_{sub}[:, d_i]}{cov(X_{sub}[:, d_i])}, \quad \forall d_i \in \{1, 2, \dots, d * w\}$

(3) $E[i, :] = -X_{normalized}[i]^2/2n, \quad \forall i$

(4) $a[i, t] = -\max(E[i, :])$

(5) $a[:, t] = \frac{a[:, t] - \min(a[:, t])}{\max(a[:, t]) - \min(a[:, t])}$

(6) $th = \min(th, \text{quantile}(a[:, t], 0.95))$

(7) $I[i, t] = 1$ if $a[i, t] > th, \quad \forall i$

As long term anomalies are of interest, time series with temporally longer anomalous behaviors are ranked more anomalous. The overall time series anomaly score A_n for each time series t_n can be computed as:

$$A_n = \frac{\sum_{t=1}^T I[n, t]}{T} \quad \forall n \quad (7)$$

For a database of time series with varying lengths, the time series anomaly score is computed by normalizing with respective lengths.

5 EXPERIMENTS

In this section, we evaluate the performance of the LAD algorithm on multi-aspect datasets. The following experiments have been conducted to study the model:

- (1) **Anomaly Detection Performance:** LAD's ability to detect real-world anomalies as compared to state-of-the-art anomaly detection models is evaluated using the ground truth labels.
- (2) **Handling Large Data:** Scalability of the LAD model on large datasets (high observation count or high dimensionality) are studied.
- (3) **Speed:** The computation and execution times of different algorithms are studied and evaluated.
- (4) **COVID-19 Time Series Data:** We study the performance of LAD model on multiple multivariate time series datasets to identify anomalous instances within each time step as well as anomalous time series amongst many.

5.1 Datasets

We consider a variety of publicly available benchmark data sets from Outlier Detection DataSets /ODDS [24] (See Tables 1) for the experimental evaluation. For the time series data, we use COVID-19 deaths and confirmed cases for US counties from John Hopkins COIVD-19 Data Repository [15].

Name	N	d	a
HTTP	567498	3	0.39%
MNIST	7603	100	9.207%
Arrhythmia	452	274	14.602%
Shuttle	49097	9	7.151%
Letter	1600	32	6.25%
Musk	3062	166	3.168%
Optdigits	5216	64	2.876%
Satellite Image	6435	36	31.639%
Speech	3686	400	1.655%
SMTp	95156	3	0.032%
Satellite Image-2	5803	36	1.224%
Forest Cover	286048	10	0.96%
KDD99	620098	29	29 0.17%

Table 1: High Dimensional and Large Sample Datasets: Description of the benchmark data sets used for evaluation of the anomaly detection capabilities of the proposed model. N - number of instances, d - number of attributes and a - fraction of known anomalies in the data set.

5.2 Baseline Methods and Parameter Initialization

As described in Section 4, LAD falls under unsupervised learning regime targeted for high dimensional data, we do not compare with supervised algorithms. For this we consider *Elliptic Envelope* (EE) [25], *Isolation Forest* (I-Forest) [19]³, *local outlier factor* (LOF) [7], and *Concentration Free Outlier Factor* CFOF [2]. The CFOF and LOF models assign an anomaly score for each data instance, while the rest of the methods provide an anomaly label. As above mentioned methods have one or more user-defined parameters, we investigated a range of values for each parameter, and report the best results. For Isolation Forest, Elliptic Envelope and CFOF, the contamination value is set to the true proportion of anomalies in the dataset.

The LAD model relies on a threshold value to classify observations with scores the value as strictly anomalous. Though this value is iteratively updated, an initial value is required by the algorithm. In this paper, the initial threshold value for the experiment is set to 0.95 for all datasets.

All the methods for anomaly detection benchmark datasets are implemented in Python and all experiments were conducted on a 2.7 GHz Quad-Core Intel Core i7 processor with a 16 GB RAM.

³The I-Forest model returns both anomaly scores and anomaly labels. As classification model outperforms its score based counterpart on above discussed datasets, we only present results on the classification model.

Table 2: Comparing LAD with existing anomaly detection algorithms for large/ high dimensional datasets using ROC-AUC as the evaluation metric.

Data	LOF	I-Forest	EE	CFOF	LAD
SHUTTLE	0.52	0.98	0.96	-	0.99
SATIMAGE-2	0.57	0.95	0.96	0.70	0.99
SATIMAGE	0.51	0.64	0.65	0.55	0.6
KDD99	0.51	0.85	0.54	-	1.0
ARRHYTHMIA	0.61	0.67	0.7	0.56	0.71
OPTDIGITS	0.51	0.52	0.45	0.49	0.48
LETTER	0.54	0.54	0.6	0.90	0.6
MUSK	0.5	0.96	0.96	0.49	0.96
HTTP	0.47	0.95	0.95	-	1.0
MNIST	0.5	0.61	0.65	0.75	0.87
COVER	0.51	0.63	0.52	-	0.96
SMTp	0.84	0.83	0.83	-	0.82
SPEECH	0.5	0.53	0.51	0.47	0.47

5.3 Evaluation Metrics

As LAD is an score based algorithm, we study the ROC curves by comparing the True Positive Rate (TPR) and False Positive Rate (FPR), across various thresholds. The final ROC-AUC (Area under the ROC curve) is reported for evaluation. For time series anomaly detection, we present the final outliers and study their deviations from normal baselines under different model settings.

5.4 Anomaly Detection Performance

Table 2 shows the performance of LOF, I-Forest, EE, CFOF and LAD on anomaly detection benchmark datasets. Due to relatively large run-time⁴, CFOF results are shown for datasets with samples less than 10k. For all the listed algorithms, results for best parameter settings are reported. The proposed LAD model outperforms other methods on most data sets. For larger and high-dimensional datasets, it can be seen from Table 2 that the LAD model outperforms all the models in most settings.⁵

To study the LAD model’s computational effectiveness, we study the computation time and scaling of LAD model on large and high dimensional datasets. We consider datasets with more than 10k observations or over 100 features for our analysis. Figures 2 and 3 show the computation time in seconds for benchmark datasets. It can be seen that the LAD model is relatively low computation time second only to Isolation Forest in most datasets. In fact, the computation time is more stable for our model as opposed to others in high dimensional datasets.

Figure 4 shows the scalability of LAD with respect to the number of records in the data. We plot the time needed to run on the first k records of the KDD-99 dataset. Each record has 29 dimensions. Figure 5 shows the scalability of LAD with respect to the number of dimensions (linear-scale). We plot the time needed to run on the first 1, 2, ..., 29 dimensions of the KDD-99 dataset. The results

⁴The CFOF model is computationally expensive relative to the rest of the algorithms. As it is aimed to study high-dimensional data, only results on datasets with <10k observations are presented.

⁵The lowest AUC values for the LAD model are observed for Speech and Optdigits data where multiple true clusters are noted.

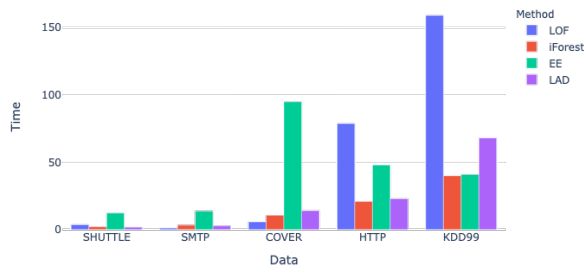


Figure 2: Computation time for large datasets

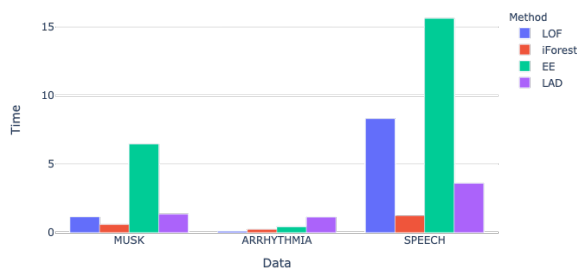


Figure 3: Computation time for high dimensional datasets

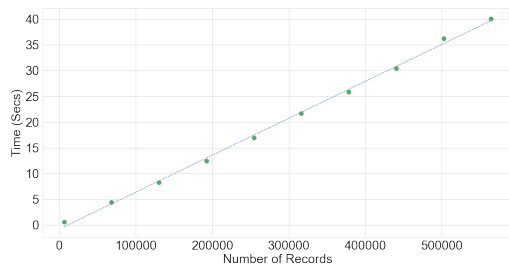


Figure 4: LAD scales linearly with the number of records for KDD-99 data

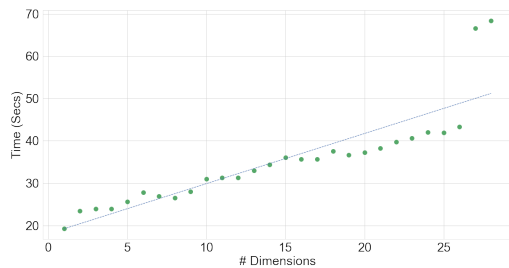


Figure 5: LAD scales linearly with the number of dimensions in KDD-99 data.

confirm the linear scalability of LAD with number of records as well as number of dimensions.

5.5 Anomaly Detection in Time Series Data

This section presents the results of LAD model on COVID-19 time series data at the US county level. Multiple settings were used to understand the data:

- (1) Deaths and confirmed case trends were considered for analysis
- (2) Daily New vs Total Counts: Both total cases as well daily new cases were analyzed for anomaly detection.
- (3) Complete history vs One Time Step: Two versions of the model were studied where data from previous time steps were and were not considered. By this, we tried to distinguish the impact of the history of the time series on identifying anomalous trends.
- (4) Univariate vs Multivariate Time Series data: To further understand the LAD model, the deaths and case trends were studied individually as a univariate time series as well as collectively in a multivariate time series data setting.
- (5) Time Series of Uniform vs Varying Lengths: Finally, all the above analyses were conducted on time series data with varying lengths. Here, for each county level time series, the time of first event was considered as initial time step to objectively study the relative temporal changes in trends.

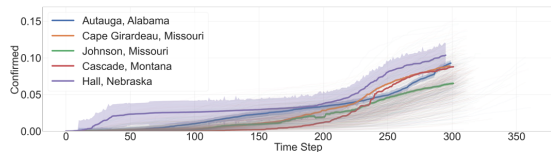
To bring all the counts to a baseline, the total counts in each time series were scaled to the respective county population. Missing information was replaced with zeros and counties with population less than 50k were eliminated from the study.

5.6 Discoveries

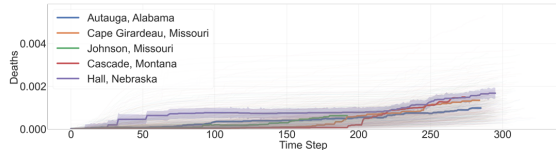
Complete history vs One Time Step. The full history setting considers the complete history of the time series and is aimed to capture most deviant trends over time. The one time step (or any smaller window) setting is more suitable to study deviations within the specific window. As we target long term deviating trends, the one time step setting returns trends that have stayed most deviant throughout the entire time range. This can be seen in Figures 6 and 7 where the one time step setting returns trends that have stayed deviant almost throughout the duration while the full history setting is able to capture significantly wider deviations. For instance, counties like Grand Forks (ND), Burleigh (ND) and Miami-Dale (FL), that had massive outbreaks at later stages⁶ were not captured as anomalous in the one time step model as seen in Figure 7a and 7b. Similarly, Hall, Nebraska, which has see a deviation in trend due to an outbreak in meat packing facility in late April 2020, was captured as anomalous trend by the full history model in Figure 6a and 6b.

Univariate vs Multivariate Time series. In Figures 6, 7, 8 and 9 we see the anomalous trends in multivariate time series, where total confirmed cases and deaths were collectively evaluated for anomaly detection. For instance, despite the near-normal trends in deaths

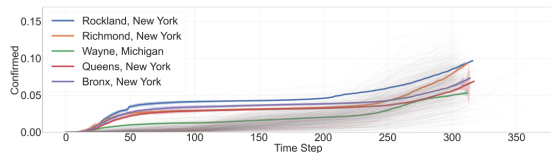
⁶<https://www.bloomberg.com/news/articles/2020-09-29/north-dakota-s-outbreak-is-as-bad-as-florida-arizona-in-july>



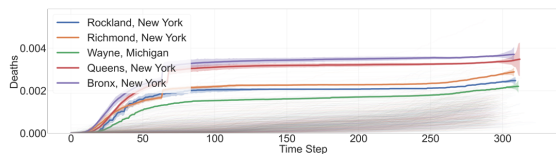
(a) Total Confirmed, Full History



(b) Total Deaths, Full History



(c) Total Confirmed, One Time Step



(d) Total Deaths, One Time Step

Figure 6: Top 5 Counties with Anomalous Trends : Varying lengths, Total Counts, Multivariate Time Series

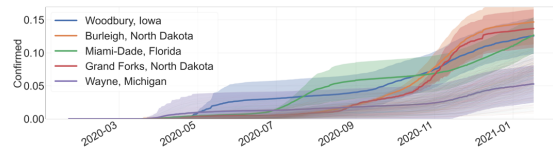
cases, Hall (NE)⁷ in Figures 6a- 6b, and Randal (TX) in Figures 9a- 9b were identified anomalous due to their the deviant confirmed case trends which significantly contributed to the anomaly scores. This setting enables identification of time-series with at least one deviating feature.

Similarly, in Figures 6c and 6d, Wayne, Michigan along with Rockland, Richmond, Queens and Bronx in NY have been identified as anomalous. In particular, Michigan was seen to have 3rd highest deaths after NY and NJ in the early stages of the pandemic with Detroit metro-area contributing to most cases⁸. Though Wayne county has near normal trend in total confirmed cases where as the total deaths trend has deviated significantly.

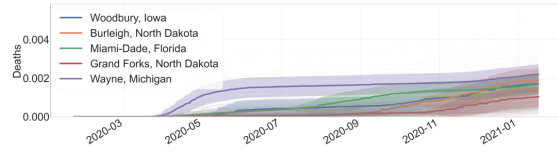
Daily New vs Total Counts. Figures 7 and 9, show anomalous trends in multivariate time series for total and daily new counts respectively. It can be seen that the anomaly score is erratic for multivariate time series on new case counts. This is due to the fact that the data for new case and death counts is more erratic leading to fluctuating normal average as well as non-smooth anomaly scores.

⁷https://www.omaha.com/news/state_and_regional/237-coronavirus-cases-tied-to-jbs-beef-plant-in-grand-island-disease-specialists-are-touring/article_2894db56-913a-5c61-a065-6860a8ae50ad.html

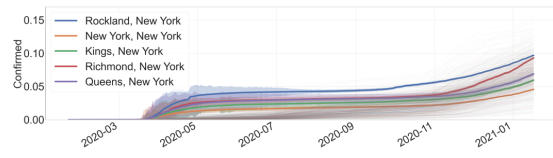
⁸<https://www.npr.org/sections/coronavirus-live-updates/2020/03/31/824738996/after-surge-in-cases-michigan-now-3rd-in-country-for-coronavirus-deaths>



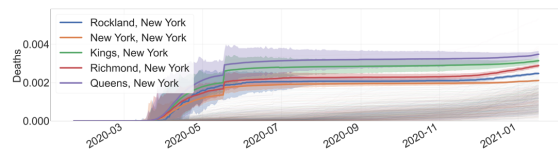
(a) Total Confirmed, Full History



(b) Total Deaths, Full History

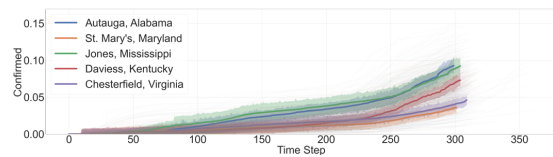


(c) Total Confirmed, One Time Step

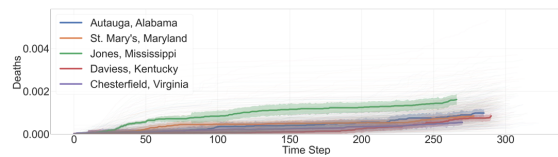


(d) Total Deaths, One Time Step

Figure 7: Top 5 Counties with Anomalous Trends : Uniform lengths, Total Counts, Multivariate Time Series



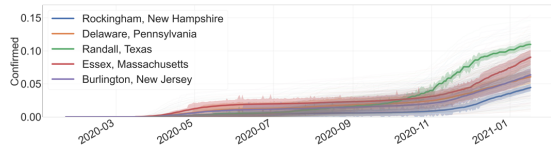
(a) New Confirmed, Full History



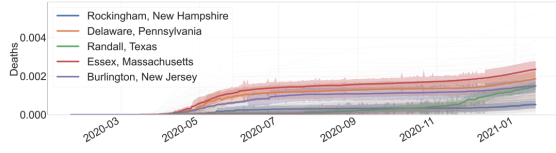
(b) New Deaths, Full History

Figure 8: Top 5 Counties with Anomalous Trends : Varying lengths, Daily New Counts, Multivariate Time Series

The LAD model on the daily new counts data was able to capture the escalation in Greater Boston area, Essex, Massachusetts in Figure 9a and 9b during March 2020. Though the total trends seem to be normal, the multiple anomalous daily trends led to their high anomaly scores. Similar patterns led to identification of Lincoln

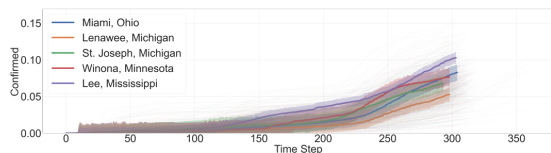


(a) New Confirmed, Full History

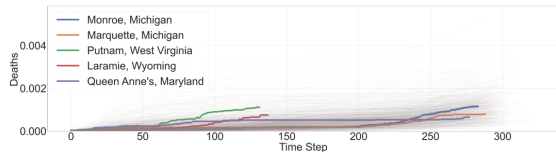


(b) New Deaths, Full History

Figure 9: Top 5 Counties with Anomalous Trends : Uniform lengths, Daily New Counts, Multivariate Time Series

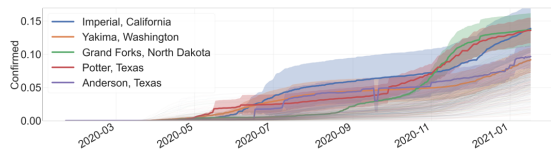


(a) Total Confirmed, Full History

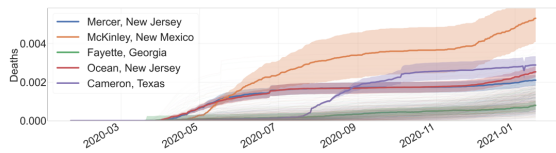


(b) Total Deaths, Full History

Figure 10: Top 5 Counties with Anomalous Trends: Varying lengths, Total counts

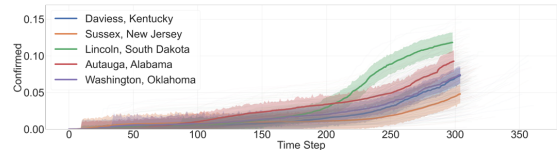


(a) Total Confirmed, Full History

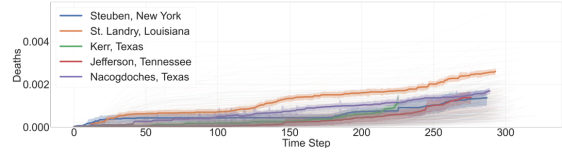


(b) Total Deaths, Full History

Figure 11: Top 5 Counties with Anomalous Trends: Uniform lengths, Total counts

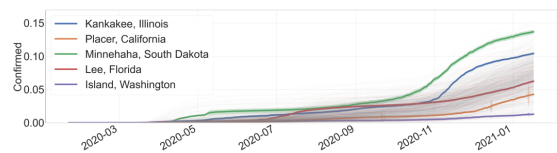


(a) New Confirmed, Full History

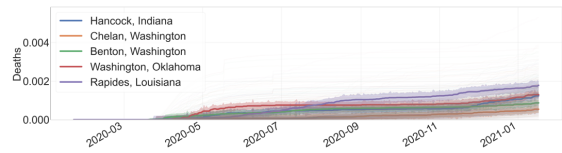


(b) New Deaths, Full History

Figure 12: Top 5 Counties with Anomalous Trends: Varying lengths, Daily New Counts



(a) New Confirmed, Full History



(b) New Deaths, Full History

Figure 13: Top 5 Counties with Anomalous Trends: Uniform lengths, Daily New Counts

(SD) and Minnehaha (SD) in Figures 12a and 13a respectively where a subsequent spike occurred after August 2020⁹.

Uniform Length vs Varying Length Time Series. The US county cases and deaths data consists of time series of uniform lengths. However, not all counties have events recorded in the early stages. Thus, studying the non-synchronized database creates a bias against counties with early reported cases. This can be seen in Figures 6 where counties like Wayne, Michigan are flagged anomalous despite starting after many counties in NY and NJ unlike in Figures 7 which reports counties in NY with an early start¹⁰. Similarly, Putnam (WV) and Laramie (WY) are found anomalous in Figure 10b where the recently evolved death trends show signs of significant divergence. On the other hand, Potter (TX) and Anderson (TX) have been identified anomalous in Figures 11a due to early increase in June 2020.

⁹<https://www.usatoday.com/story/news/nation/2020/08/07/sturgis-motorcycle-rally-what-know-masks-attendance-rules/3321223001/>

¹⁰<https://www.npr.org/sections/coronavirus-live-updates/2020/03/31/824738996/after-surge-in-cases-michigan-now-3rd-in-country-for-coronavirus-deaths>

6 CONCLUSION

In this paper, we propose LAD, a novel scoring algorithm for anomaly detection in large/high-dimensional data. The algorithm successfully handles high dimensions by implementing large deviation theory. Our contributions include reestablishing the advantages of large deviations theory to large and high dimensional datasets. We also present an online extension of the model that is aimed to identify anomalous time series in a multivariate time series data. The model shows vast potential in scalability and performance against baseline methods. The online LAD returns a temporally evolving score for each time series that allows us to study the deviations in trends relative to the complete time series database.

A potential extension to the model could include anomalous event detection for each individual time series. Another possible future work could be extending the model to enable anomaly detection in multi-modal datasets. Additionally, the online LAD model could be enhanced to use temporally weighted scores prioritizing recent events.

REFERENCES

- [1] Charu C Aggarwal and Philip S Yu. 2001. Outlier detection for high dimensional data. In *Proceedings of the 2001 ACM SIGMOD international conference on Management of data*. 37–46.
- [2] Fabrizio Angiulli. 2020. CFOF: a concentration free measure for anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 14, 1 (2020), 1–53.
- [3] Fabrizio Angiulli and Clara Pizzuti. 2002. Fast outlier detection in high dimensional spaces. In *European conference on principles of data mining and knowledge discovery*. Springer, 15–27.
- [4] Laura Beggel, Bernhard X Kausler, Martin Schiegg, Michael Pfeiffer, and Bernd Bischl. 2019. Time series anomaly detection based on shapelet learning. *Computational Statistics* 34, 3 (2019), 945–976.
- [5] Seif-Eddine Benkabou, Khalid Benabdeslem, and Bruno Canitia. 2018. Unsupervised outlier detection for time series by entropy and dynamic time warping. *Knowledge and Information Systems* 54, 2 (2018), 463–486.
- [6] Carl Boettiger and Alan Hastings. 2013. No early warning signals for stochastic transitions: insights from large deviation theory. *Proceedings of the Royal Society B: Biological Sciences* 280, 1766 (2013), 20131372.
- [7] Markus M. Breunig, Hans-Peter Kriegel, Raymond T. Ng, and J. Sander. 2000. LOF: identifying density-based local outliers. In *Proceedings of 2000 ACM SIGMOD International Conference on Management of Data*. 93–104.
- [8] Zeynep Ceylan. 2020. Estimation of COVID-19 prevalence in Italy, Spain, and France. *Science of The Total Environment* 729 (2020), 138817.
- [9] Varun Chandola, Arindam Banerjee, and Vipin Kumar. 2009. Anomaly detection: A survey. *Comput. Surveys* 41, 3 (2009).
- [10] Varun Chandola, Deepthi Cheboli, and Vipin Kumar. 2009. *Detecting Anomalies in a Timeseries Database*. Technical Report 09-004. University of Minnesota, Computer Science Department.
- [11] V. Chandola and V. Kumar. 2008. A Comparative Evaluation of Anomaly Detection Techniques for Sequence Data. In *Proceedings of International Conference on Data Mining*. Pisa, Italy.
- [12] Sanjay Chawla and Aristides Gionis. 2013. k-means-: A unified approach to clustering and outlier detection. In *SDM*.
- [13] Giovanni Dematteis, Tobias Grafke, and Eric Vanden-Eijnden. 2018. Rogue waves and large deviations in deep sea. *Proceedings of the National Academy of Sciences* 115, 5 (2018), 855–860.
- [14] Frank Den Hollander. 2008. *Large deviations*. Vol. 14. American Mathematical Soc.
- [15] Ensheng Dong, Hongru Du, and Lauren Gardner. 2020. An interactive web-based dashboard to track COVID-19 in real time. *The Lancet infectious diseases* 20, 5 (2020), 533–534.
- [16] Tak-chung Fu. 2011. A review on time series data mining. *Engineering Applications of Artificial Intelligence* 24, 1 (2011), 164–181.
- [17] Manish Gupta, Jing Gao, Charu C Aggarwal, and Jiawei Han. 2013. Outlier detection for temporal data: A survey. *IEEE Transactions on Knowledge and data Engineering* 26, 9 (2013), 2250–2267.
- [18] Victoria Hodge and Jim Austin. 2004. A Survey of Outlier Detection Methodologies. *Artificial Intelligence Review* 22, 2 (2004), 85–126.
- [19] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2012. Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 6, 1 (2012), 1–39.
- [20] Mohsen Maleki, Mohammad Reza Mahmoudi, Darren Wraith, and Kim-Hung Pho. 2020. Time series modelling to forecast the confirmed and recovered cases of COVID-19. *Travel medicine and infectious disease* 37 (2020), 101742.
- [21] Thomas Mikosch and Olivier Wintenberger. 2016. A large deviations approach to limit theory for heavy-tailed time series. *Probability Theory and Related Fields* 166, 1 (2016), 233–269.
- [22] Ioannis Ch Paschalidis and Georgios Smaragdakis. 2008. Spatio-temporal network anomaly detection by assessing deviations of empirical measures. *IEEE/ACM Transactions On Networking* 17, 3 (2008), 685–697.
- [23] Sridhar Ramaswamy, Rajeev Rastogi, and Kyuseok Shim. 2000. Efficient algorithms for mining outliers from large data sets. In *Proceedings of the 2000 ACM SIGMOD international conference on Management of data*. ACM Press, Dallas, Texas, United States, 427–438.
- [24] Shebuti Rayana. 2016. ODDS Library. <http://odds.cs.stonybrook.edu>
- [25] Peter J Rousseeuw and Katrien Van Driessen. 1999. A fast algorithm for the minimum covariance determinant estimator. *Technometrics* 41, 3 (1999), 212–223.
- [26] Hugo Touchette. 2009. The large deviation approach to statistical mechanics. *Physics Reports* 478, 1-3 (2009), 1–69.
- [27] SR Srinivasa Varadhan. 1984. *Large deviations and applications*. SIAM.
- [28] SR Srinivasa Varadhan. 2010. Large deviations. In *Proceedings of the International Congress of Mathematicians 2010 (ICM 2010) (In 4 Volumes) Vol. I: Plenary Lectures and Ceremonies Vols. II–IV: Invited Lectures*. World Scientific, 622–639.
- [29] Dragomir Yankov, Eamonn Keogh, and Umaa Rebbapragada. 2008. Disk aware discord discovery: Finding unusual time series in terabyte sized datasets. *Knowledge and Information Systems* 17, 2 (2008), 241–262.
- [30] Abdelhafid Zeroual, Fouzi Harrou, Abdelkader Dairi, and Ying Sun. 2020. Deep learning methods for forecasting COVID-19 time-Series data: A Comparative study. *Chaos, Solitons & Fractals* 140 (2020), 110121.
- [31] Chuxu Zhang, Dongjin Song, Yuncong Chen, Xinyang Feng, Cristian Lumezanu, Wei Cheng, Jingchao Ni, Bo Zong, Haifeng Chen, and Nitesh V Chawla. 2019. A deep neural network for unsupervised anomaly detection and diagnosis in multivariate time series data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 1409–1416.