
Anomaly detection on streamed data

Thomas Cochrane

The Alan Turing Institute, London
thomasc@turing.ac.uk

Peter Foster

The Alan Turing Institute, London
pfoster@turing.ac.uk

Terry Lyons

Mathematical Institute, University of Oxford
The Alan Turing Institute, London
terry.lyons@maths.ox.ac.uk

Imanol Perez Arribas

Mathematical Institute, University of Oxford
The Alan Turing Institute, London
Oxford, United Kingdom
imanol.perez@maths.ox.ac.uk

Abstract

We introduce powerful but simple methodology for identifying anomalous observations against a corpus of ‘normal’ observations. All data are observed through a vector-valued feature map. Our approach depends on the choice of corpus and that feature map but is invariant to affine transformations of the map and has no other external dependencies, such as choices of metric; we call it conformance. Applying this method to (signatures) of time series and other types of streamed data we provide an effective methodology of broad applicability for identifying anomalous complex multimodal sequential data. We demonstrate the applicability and effectiveness of our method by evaluating it against multiple data sets. Based on quantifying performance using the receiver operating characteristic (ROC) area under the curve (AUC), our method yields an AUC score of 98.9% for the PenDigits data set; in a subsequent experiment involving marine vessel traffic data our approach yields an AUC score of 89.1%. Based on comparison involving univariate time series from the UEA & UCR time series repository with performance quantified using balanced accuracy and assuming an optimal operating point, our approach outperforms a state-of-the-art shapelet method for 19 out of 28 data sets.

1 Introduction

The task of *anomaly detection*, i.e. the task of determining whether a given observation is *unusual* compared to a corpus of observations deemed to be *normal* or *usual*, is a challenge with applications in various fields such as medicine [10], financial fraud [21] and cybersecurity [14].

The idea of using a metric to discriminate a corpus of scenarios from anomalies, and the view that an event is an anomaly if it is some distance from the set of observations, seem natural; these have been used many times [7]. The main weakness of this approach is the arbitrariness of the choice of metric and, possibly, some way of calibrating the power of the technique. An important innovation in this paper is the use of the variance, the dual norm to the covariance, as the metric. As we will explain, in many way it is surprising that the choice works, but in fact there is a strong and quite deep mathematical explanation for its effectiveness in terms of concentration of measure. It is a measure of exceptionality that can be applied to any corpus of data described through a vector feature set. It also provides internal measures of its own effectiveness in terms of the extent to which members of the corpus are themselves anomalies to the rest of the corpus. It requires no external choices or parameters. For example, linear transformations of the features do not change the analysis or the measures of exceptionality at all.

1.1 Existing work

Anomaly detection comprises a vast literature spanning multiple disciplines [7]. Among unsupervised anomaly detection techniques applicable to multivariate data, existing work includes density-based approaches [5], clustering [11], random forests [18], support vector machines [2] and neural networks [6].

The time series anomaly detection literature has largely focused on detecting anomalous points or subsequences within a time series, rather than detecting entire time series as anomalous. Hyndman [13] detects anomalous time series by calculating a set of features of the overall time series, and projecting to the two principal components; Beggel et al. [4] learn shapelet-based features that are particularly associated with the normal class.

1.2 Our work

There are many data science contexts where it is already meaningful to construct vector representations or features to describe data. Word2Vec [20] and Kernels [12] provide two examples. The method introduced here could easily be applied to these contexts. In this paper, we initially, and specifically, focus on the signature as a vectorisation for streamed data, establishing that the methods are easy to apply and effective.

Definition 1.1 (Variance norm). Let μ be a probability measure on a vector space V . The covariance quadratic form $\text{Cov}(\psi, \phi) := \mathbb{E}^\mu[\psi(x)\phi(x)]$, defined on the dual of V , induces a dual norm defined for $x \in V$ by

$$\|x\|_\mu := \sup_{\text{Cov}(\phi, \phi) \leq 1} \phi(x) \quad (1)$$

on $x \in V$. It is finite on the linear span of the support of μ , and infinite outside of it. We refer to this norm, computed for the measure μ re-centered to have mean zero, as the variance norm $\|\cdot\|_\mu$ associated to μ .

The variance norm is well defined whenever the measure has finite second moments and, in particular, for the empirical measure associated to a finite set of observations $\{x_i\}_i$.

This variance is surprisingly useful for detecting anomalies. Consider the standard normal distribution in d dimensions. That is to say consider $Z := (z_1, \dots, z_d)$ where the z_i are independent normal variables with mean zero and variance one. Then the covariance is the usual Euclidean inner product, and the variance of Z is the usual Euclidean norm. Note that the expected value of $x \in V$ by $\|x\|_\mu^2$ is d , and as $d \rightarrow \infty$ the norm of the path is of the order of \sqrt{N} , converging to infinity with N . In the high dimensional case, the norm is huge and in the infinite dimension case it is infinite. For Brownian motion on $[0, T]$ it is the L^2 norm of the gradient of the path. Clearly no Brownian path has a derivative, let alone a square integrable one. We see that the variance is intrinsic, but actually provides a very demanding notion of nearness that looks totally unrealistic.

However, keeping the Gaussian context, there is a fundamental theorem in stochastic analysis known as the TSB isoperimetric inequality¹ [1]. An immediate corollary is that if one takes any set A of probability one half in V , and a new sample Z from the Gaussian measure, then the probability that Z is a variance-distance r from A is at most $1/\sqrt{2\pi} \int_r^\infty \exp(u^2/2t) du$ and so vanishing small if r is even of moderate size. A Brownian path may be irregular, but if you take a corpus of Brownian paths with probability at least a half, then it will differ from one of those paths by a differentiable path of small norm. This makes the variance an excellent measure of exceptionality, it is selective and discriminatory, but it must be used to compare with the corpus and not used directly. A new member Z of the corpus will be far away from most members of the corpus, but there will with very high probability be some members of the corpus to which it is associated very well. With this in mind we make the following definition:

Definition 1.2. Let μ be a probability measure on a vector space V . Define the *conformance* of x to μ to be the distance

$$\text{dist}(x; \mu) := \inf_{y \in \text{supp}(\mu)} \|x - y\|_\mu.$$

¹TSB stands for Tsirelson-Sudakov-Borell.

If $S : V \rightarrow W$ is a linear map, then $S(x) \in W$ is more conformant to $S(\mu)$ than $x \in V$ is to μ (the conformance score is reduced by the linear map).

Keeping to the Gaussian context, let A be any set of measure $1/2$ and let μ be the Gaussian measure restricted to A and normalised. Then, reiterating, the TSB inequality ensures conformance to μ is an excellent measure of exceptionality.

An empirical measure is not in itself Gaussian, even if drawn from a Gaussian. So taking half of the ensemble only captures the other half tightly when the sample size is large enough compared with the dimension of the feature set that balls round it capture a good proportion of the probability measure. Before that, the resolution provided by the feature set is so high that essentially every sample is non-conformant. Fortunately this is easy to measure empirically by looking to identify R . Therefore, if we split the corpus randomly into two halves, the probability is one half that a point chosen from the second half of the corpus is within a distance R of the first half. From that scale on, if the conformance of a new observation to the corpus is $r + R$, then r should provide an effective measure of being an anomaly and R provides a measure of the extent to which the dimension of the feature set is overwhelming the sample size and in this context every observation is innovative and an anomaly.

The non-Gaussian case lacks the very sharp theoretical underpinning of the Gaussian case, but the approach remains clear and its power can still easily be determined from the data. We validate the approach by identifying anomalies in streamed data using signatures as the vector features.

Our methodology provides a data-driven notion of a *distance* (i.e. conformance) between an arbitrary stream of data and the corpus. Moreover, it has four properties that are particularly useful for anomaly detection:

- The variance norm is intrinsic to the vector representation and independent of any choice of basis.
- The conformance score, as a measure of anomaly, does not depend on any external choice of metric, etc.
- By using the signature to vectorise the corpus of streamed data, it is straightforward to accommodate streams that are differently sampled and essentially multimodal.
- There are no distribution assumptions on the corpus of vectors.

The paper is structured as follows: Section 2 introduces the basic signature tools. In Section 3 we combine conformance with signatures to analyse anomalies in streamed data. In Section 4, we report results of numerical experiments on PenDigits, marine vessel traffic data and univariate time series from the UEA & UCR repository. In Section 5 we briefly summarise the contribution of the paper.

2 Streams of data and signature features

2.1 Streams of data

Below we give a formal definition of a stream of data, [16, Definition 2.1].

Definition 2.1 (Stream of data). The space of streams of data in a set \mathcal{X} is defined as

$$\mathcal{S}(\mathcal{X}) := \{\mathbf{x} = (x_1, \dots, x_n) : x_i \in \mathcal{X}, n \in \mathbb{N}\}.$$

Example 2.2. When a person writes a character by hand, the stroke of the pen naturally determines a path. If we record the trajectory we obtain a two-dimensional stream of data $\mathbf{x} = ((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \in \mathcal{S}(\mathbb{R}^2)$. If we record the stroke of a different writer, the associated stream of data could have a different number of points. The distance between successive points may also vary.

2.2 Signature features

Definition 2.3 (Signature). Let $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{S}(\mathbb{R}^d)$ be a stream of data in d dimensions. Let $X = (X_1, \dots, X_d) : [0, 1] \rightarrow \mathbb{R}^d$ be such that $X\left(\frac{i}{n-1}\right) = x_{i+1}$ for $i = 0, 1, \dots, n-1$ and linear

interpolation in between. Then, we define the signature of \mathbf{x} of order $N \in \mathbb{N}$ as

$$\text{Sig}^N(\mathbf{x}) := \left(\int_{0 < t_1 < \dots < t_k < 1} \frac{dX_{i_1}}{dt}(t_1) \cdot \frac{dX_{i_2}}{dt}(t_2) \cdots \frac{dX_{i_k}}{dt}(t_k) dt_1 \cdots dt_k \right)_{\substack{1 \leq i_1, \dots, i_k \leq d \\ k=0,1,2,\dots,N}}. \quad (2)$$

The signature of a stream of data is a vector of scalars. The dimension of this vector is

$$d_N := 1 + d + d^2 + \dots + d^N = \frac{d^{N+1} - 1}{d - 1}.$$

Proposition 2.4. For each $N \in \mathbb{N}$, define $d_N := \frac{d^{N+1} - 1}{d - 1}$, which is the dimension of the signature of order N . There exists a product

$$\mathfrak{w} : \mathbb{R}^{d_N} \times \mathbb{R}^{d_N} \rightarrow \mathbb{R}^{d_{2N}}$$

called the shuffle product such that

$$\langle f, \text{Sig}^N(\mathbf{x}) \rangle \langle g, \text{Sig}^N(\mathbf{x}) \rangle = \langle f \mathfrak{w} g, \text{Sig}^{2N}(\mathbf{x}) \rangle \quad \forall f, g \in \mathbb{R}^{d_N}, \mathbf{x} \in \mathcal{S}(\mathbb{R}^d),$$

where $\langle \cdot, \cdot \rangle$ denotes the inner (dot) product.

See [15, Definition 2.5] for an explicit construction of the shuffle product \mathfrak{w} .

2.3 Stream transformations

Stream transformations map a stream of data to another stream of data that one considers might contain relevant information for the problem at hand.

Definition 2.5. A stream transformation is a mapping

$$T : \mathcal{S}(\mathbb{R}^d) \rightarrow \mathcal{S}(\mathbb{R}^{d'}),$$

where typically $d' \geq d$.

Below we introduce a few stream transformations that have proved to be popular with signatures in the literature, and will be used in later sections. More than one transformation can simultaneously be applied on a single stream.

2.3.1 Time transformation

The time transformation adds an extra dimension to a stream of data, which accounts for time:

$$\begin{aligned} T_{\text{time}} : \mathcal{S}(\mathbb{R}^d) &\rightarrow \mathcal{S}(\mathbb{R}^{1+d}) \\ (x_0, \dots, x_i, \dots, x_n) &\mapsto ((t_0, x_0), \dots, (t_i, x_i), \dots, (t_n, x_n)). \end{aligned}$$

where t_0, \dots, t_n are chosen to be strictly increasing. Assume our data includes timestamps t_0, \dots, t_n . A variant of this transformation involves computing differences between successive timestamps:

$$\begin{aligned} T_{\text{time-diff}} : \mathcal{S}(\mathbb{R}^d) &\rightarrow \mathcal{S}(\mathbb{R}^{1+d}) \\ (x_0, \dots, x_i, \dots, x_n) &\mapsto ((0, x_0), \dots, (t_i - t_{i-1}, x_i), \dots, (t_n - t_{n-1}, x_n)). \end{aligned}$$

2.3.2 Lead-lag transformation

The lead-lag transformation of a d -dimensional stream of data of length n is a $2d$ -dimensional stream of data of length $2n - 1$, defined as follows:

$$\begin{aligned} T_{\text{lead-lag}} : \mathcal{S}(\mathbb{R}^d) &\rightarrow \mathcal{S}(\mathbb{R}^{2d}) \\ (x_0, \dots, x_i, \dots, x_n) &\mapsto (\hat{x}_0, \dots, \hat{x}_i, \dots, \hat{x}_{2n}) \end{aligned}$$

where

$$\hat{x}_{2i} = (x_i, x_i) \in \mathbb{R}^{2d}, \quad \hat{x}_{2i+1} = (x_i, x_{i+1}) \in \mathbb{R}^{2d}$$

for $i = 0, 1, \dots, n$. The work in [9] studies the signature of lead-lag transformed streams.

2.3.3 Invisibility transform

Signatures are constructed from increments of the stream of data. As a consequence, all information about the absolute value of the steps of the stream is lost. Sometimes it is desirable to keep reference to the absolute value of the underlying stream; in this case the *invisibility transform* [23] is useful. When taking the signature of a stream after applying the invisibility transform, the absolute value of the stream is preserved.

The invisibility transform is defined as follows:

$$\begin{aligned} T_{inv} : \mathcal{S}(\mathbb{R}^d) &\rightarrow \mathcal{S}(\mathbb{R}^{d+1}) \\ (x_0, \dots, x_i, x_n) &\mapsto (\hat{x}_0, \dots, \hat{x}_i, \dots, \hat{x}_{n+1}) \end{aligned}$$

where

$$\hat{x}_0 = (x_0, 0) \in \mathbb{R}^{d+1}, \quad \hat{x}_i = (x_{i-1}, 1) \in \mathbb{R}^{d+1}$$

for $i = 1, \dots, n + 1$.

3 Anomalies in streamed data

Let $\mathcal{C} \subset \mathcal{S}(\mathbb{R}^d)$ be a finite corpus (or empirical measure) of streams of data. Let Sig^N be the signature of order $N \in \mathbb{N}$. Then $\|\cdot\|_{\text{Sig}^N(\mathcal{C})}$ is the variance norm associated with the empirical measure of $\text{Sig}^N(\mathcal{C})$.

There are a number of interesting relationships between the variance norm and the signature, one being ease of computation; the variance norm of order N can easily be computed from the expected signature of order $2N$.

Proposition 3.1. *Let $w \in \mathbb{R}^{d_N}$. We have*

$$\|w\|_{\text{Sig}^N(\mathcal{C})} = \sup_{f \in \mathbb{R}^{d_N} \setminus \{0\}} \frac{\langle f, w \rangle^2}{\langle f, w \rangle^2, \mathbb{E}[\text{Sig}^{2N}(\mathbf{x})]}.$$

Proposition 3.2. *Let $A_{i,j} := \langle e_i \omega e_j, \mathbb{E}[\text{Sig}^{2N}(\mathbf{x})] \rangle$ for $i, j = 1, \dots, d_N$. Then,*

$$\|w\|_{\text{Sig}^N(\mathcal{C})} = \langle w, A^{-1}w \rangle \quad \forall w \in \mathbb{R}^{d_N}.$$

Appendix C describes some other interesting properties.

3.1 Anomaly detection using conformance

Let $\mathcal{D} \subset V$ be a finite corpus of vector data. Use a large conformance score (Definition 1.2) to identify outlying behaviour. As explained above, each corpus has its own threshold of conformance. So, we randomly split the corpus into two equal-sized parts and denote the empirical probability measures on those two parts by \mathcal{D}_1 and \mathcal{D}_2 . For a random point $x \in V$ with law \mathcal{D}_2 we can look at its conformance to \mathcal{D}_1 . By looking at the right \mathcal{D}_2 -tail of the random variable $\text{dist}(x; \mathcal{D}_1)$ with a given probability $\varepsilon > 0$ we have a natural quantified choice of anomalous behaviour. A point chosen randomly from \mathcal{D}_2 has a probability of at most ε of a conformance that exceeds the threshold.

Depending on the choice of vector feature map for the corpus the power of this approach will change. For example, if the feature map is very high-dimensional, the threshold will have poor discriminatory power. The same is true for very low-dimensional feature maps. This is where, in the context of streamed data, the graded nature of the signature features proves to be advantageous.

4 Evaluation

We apply our method to the task of unsupervised anomaly detection. That is, we have a data set $\mathcal{I} \subset \mathcal{S}(\mathbb{R}^n)$ partitioned into those data deemed to be normal $\mathcal{I}_{\text{normal}}$ and those data deemed to be anomalous $\mathcal{I}_{\text{anomaly}}$. By further partitioning, we obtain the corpus $\mathcal{C} \subset \mathcal{I}$ which we use for training; as our testing data \mathcal{Y} we use $\mathcal{Y} := \mathcal{I} \setminus \mathcal{C}$.

We perform experiments on a 2018 MacBook Pro equipped with a 2.6 GHz 6-Core Intel Core i7 processor and 32 GB 2400 MHz DDR4 memory. For the results reported in Table 1, Table 2, Figure 2, the respective CPU times observed are 54min, 2d 3h 51min, 4h 59min. To compute signatures of streams, we use the `iisignature` library [22].

4.1 Handwritten digits

We evaluate our proposed method using the PenDigits-orig data set [8]. This data set consists of 10992 instances of hand-written digits captured from 44 subjects using a digital tablet and stylus, with each digit represented approximately equally frequently. Each instance is represented as a 2-dimensional stream, based on sampling the stylus position at 10Hz.

We apply the PenDigits data to unsupervised anomaly detection by defining $\mathcal{I}_{\text{normal}}$ as the set of instances representing digit m . We define \mathcal{C} as the subset of $\mathcal{I}_{\text{normal}}$ labelled as ‘training’ by the annotators. Furthermore, we define \mathcal{Y} as the set of instances labelled as ‘testing’ by the annotators ($|\mathcal{Y}| = 3498$). Finally, we define $\mathcal{I}_{\text{anomaly}}$ as the subset of \mathcal{Y} not representing digit m . Considering all possible digits, we obtain on average $|\mathcal{C}| = 749.4$, $|\mathcal{I}_{\text{anomaly}}| = 3148.2$. Assuming that digit class is invariant to translation and scaling, we apply Min-Max normalisation to each individual stream.

Table 1 displays results based on taking signatures of order $N \in [1..5]$ and without any stream transformations applied. The results are based on aggregating conformance values across the set of possible digits before computing the ROC AUC. As we observe, performance increases monotonically from 0.901 ($N = 1$) to 0.989 ($N = 5$). Figure 3 displays plots of empirical cumulative distributions of conformance values that we obtain for normal and anomalous testing data across values of N .

$N = 1$	$N = 2$	$N = 3$	$N = 4$	$N = 5$
0.901 \pm 0.004	0.965 \pm 0.002	0.983 \pm 0.001	0.987 \pm 0.001	0.989 \pm 0.000

Table 1: Handwritten digits data: performance quantified using ROC AUC in response to signature order N . Confidence intervals are bootstrapped standard errors based on 10^5 samples.

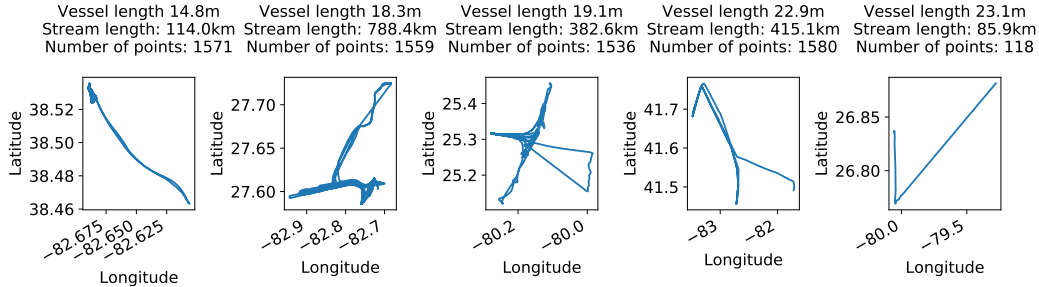
4.2 Marine vessel traffic data

Next, we consider a sample of marine vessel traffic data², based on the automatic identification system (AIS) which reports a ship’s geographical position alongside other vessel information. The AIS data that we consider were collected by the US Coast Guard in January 2017, with a total of 31 884 021 geographical positions recorded for 6 282 distinct vessel identifiers. We consider the stream of timestamped latitude/longitude position data associated with each vessel a representation of the vessel’s path. Figure 1 displays stream data for a sample of vessels.

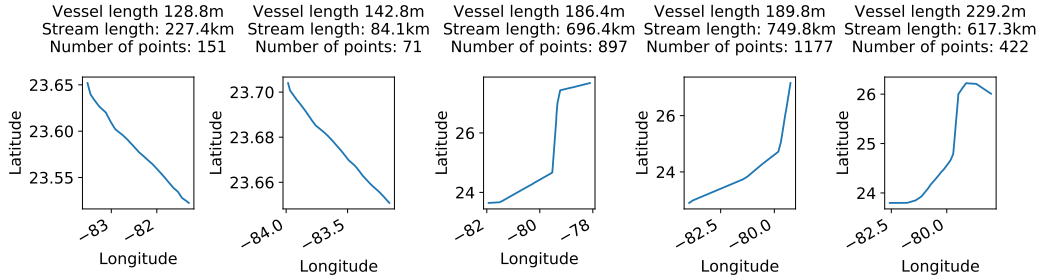
We prepare the marine vessel data by retaining only those data points with a valid associated vessel identifier. In addition, we discard vessels with any missing or invalid vessel length information. Next, to help constrain computation time, we compress each stream by retaining a given position only if its distance relative to the previously retained position exceeds a threshold of 10m. Finally, to help ensure that streams are faithful representations of ship movement, we retain only those vessels whose distance between initial and final positions exceeds 5km. To evaluate the effect of stream length on performance, we disintegrate streams so that the length D between initial and final points in each sub-stream remains constant with $D \in \{4\text{km}, 8\text{km}, 16\text{km}, 32\text{km}\}$. After disintegrating streams, we retain only those sub-streams whose maximum distance between successive points is less than 1km.

We partition the data by deeming a sub-stream normal if it belongs to a vessel with a reported vessel length greater than 100m. Conversely, we deem sub-streams anomalous if they belong to vessels with a reported length less than or equal to 50m. We obtain the corpus \mathcal{C} from 607 vessels, whose sub-streams total between 10 111 ($D = 32\text{km}$) and 104 369 ($D = 4\text{km}$); we obtain the subset of normal instances used for testing $\mathcal{I}_{\text{normal}} \setminus \mathcal{C}$ from 607 vessels, whose sub-streams total between 11 254 ($D = 32\text{km}$) and 114 071 ($D = 4\text{km}$); lastly we obtain the set of anomalous instances $\mathcal{I}_{\text{anomaly}}$ from 997 vessels whose sub-streams total between 8 890 ($D = 32\text{km}$) and 123 237 ($D = 4\text{km}$). To

²https://coast.noaa.gov/htdata/CMSP/AISDataHandler/2017/AIS_2017_01_Zone17.zip, accessed May 2020.



(a) Vessels with reported length less than 50m



(b) Vessels with reported length greater than 100m

Figure 1: Sample of marine vessel paths.

account for any imbalance in the number of sub-streams associated with vessels, we use for each of the aforementioned three subsets a weighted sample of 5 000 instances.

After computing sub-streams and transforming them as described in Sections 2.3.2 and 2.3.3, we apply Min-Max normalisation with respect to the corpus \mathcal{C} . To account for velocity, we incorporate the difference between successive timestamps as an additional dimension, as described in Section 2.3.1.

We report results based on taking signatures of order $N = 3$. For comparison, as a baseline approach we summarise each sub-stream by estimating its component-wise mean and covariance, retaining the upper triangular part of the covariance matrix. This results in feature vectors of dimensionality $\frac{1}{2}(n^2 + 3n)$ which we provide as input to an isolation forest [17]. We train the isolation forest using 100 trees and for each tree in the ensemble using 256 samples represented by a single random feature.

Table 2 displays results for our proposed approach in comparison to the baseline, for combinations of stream transformations and values of the sub-stream length D . Signature conformance yields higher ROC AUC scores than the baseline for 30 out of 32 parameter combinations. The maximum ROC AUC score of 0.891 is for a combination of lead-lag, time differences, and invisibility reset transformations with $D = 32\text{km}$, using the signature conformance. Compared to the best-performing baseline parameter combination, this represents a performance gain of 6.8 percentage points.

4.3 Univariate time series

For the specific case of detecting anomalous univariate time series, we benchmark our method against the ADSL shapelet method of Beggel et al. [4], using their set of 28 data sets from the UEA & UCR time series repository [3] adapted in exactly the same manner. Each data set comprises a set of time series of equal length, together with class labels. One class (the same as in ADSL) is designated as a normal class, with all other classes designated as anomalies. To prepare the data for our method, we convert each time series into a 2-dimensional stream by incorporating a uniformly-increasing time dimension. We apply no other transformations to the data, and take signatures of order $N = 5$.

We create training and test sets exactly as in ADSL. The training corpus \mathcal{C} consists of 80% of the normal time series, contaminated by a proportion of anomalies (we compute results for anomaly rates of 0.1% and 5%). Across these data sets $|\mathcal{C}|$ ranges from 10 (Beef) to 840 (ChlorineConcentration at 5%), $|\mathcal{I}_{\text{normal}}|$ ranges from 2 (Beef) to 200 (ChlorineConcentration), and $|\mathcal{I}_{\text{anomaly}}|$ ranges from 19

Transformation			Conformance dist(\cdot ; $\text{Sig}^3(C)$)				Isolation forest baseline			
Lead-lag	Time-Diff	Inv. Reset	Sub-stream length D				Sub-stream length D			
			4km	8km	16km	32km	4km	8km	16km	32km
No	No	No	0.723	0.706	0.705	0.740	0.690	0.718	0.717	0.733
No	No	Yes	0.776	0.789	0.785	0.805	0.682	0.698	0.714	0.716
No	Yes	No	0.810	0.813	0.818	0.848	0.771	0.779	0.779	0.803
No	Yes	Yes	0.839	0.860	0.863	0.879	0.745	0.751	0.761	0.797
Yes	No	No	0.811	0.835	0.824	0.837	0.759	0.765	0.766	0.763
Yes	No	Yes	0.812	0.835	0.833	0.855	0.755	0.761	0.763	0.762
Yes	Yes	No	0.845	0.861	0.862	0.877	0.820	0.815	0.823	0.817
Yes	Yes	Yes	0.848	0.863	0.870	0.891	0.810	0.795	0.816	0.815

Table 2: Marine vessel traffic data: performance quantified using ROC AUC for combinations of stream transformations and sub-stream length D . For each parameter combination, bold is best between signature and baseline. Italics is global best. The bootstrapped standard errors based on 10^5 samples range between 0.003 and 0.005.

(BeetleFly and BirdChicken at 0.1%) to 6401 (Wafer at 5%). We run experiments with ten random train-test splits, and take the median result. The performance measure used by ADSL is the balanced accuracy, which requires a threshold to be set for detecting anomalies. We report the best achievable balanced accuracy across all possible thresholds, and compare against the best value reported for ADSL. Figure 2 plots our results. Individual scores are available in Table 3.

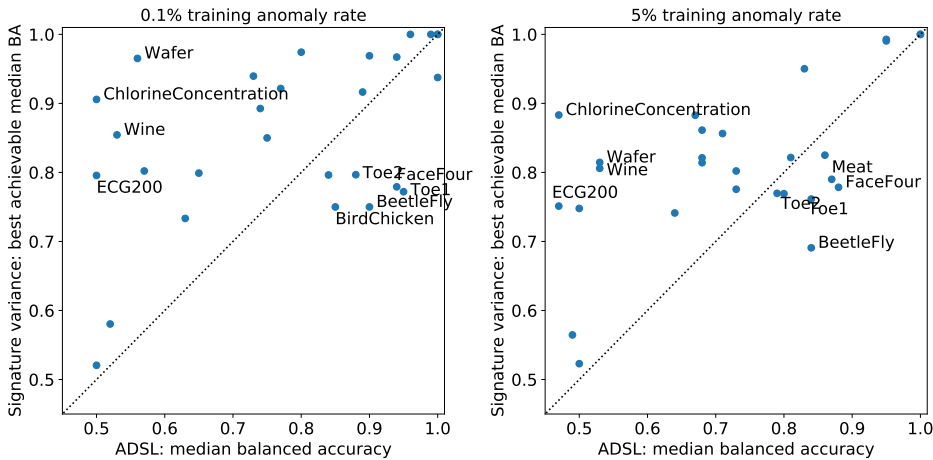


Figure 2: Comparison of our method against ADSL [4]

Our method performs competitively with ADSL, both when the proportion of anomalies in the training corpus is low and when it is high. It is able to detect anomalies in four of the six data sets where ADSL struggles because the anomalies are less visually distinguishable (ChlorineConcentration, ECG200, Wafer, Wine). However, there are data sets where ADSL performs better (BeetleFly, BirdChicken, FaceFour, ToeSegmentation1 and ToeSegmentation2): these data sets largely originate from research into shapelet methods, and they appear to contain features that are detected well by shapelets. Applying transformations to the data sets before input may improve our method’s results.

5 Conclusion

Motivated by the TSB isoparametric inequality we introduce the notion of conformance as an intrinsic and canonical tool to identify anomalous behaviour. It seems well-matched to the important challenge of identifying anomalous trajectories of streamed data against a corpus of ‘normality’. The approach appears robust when tested against a wide variety of data sets.

The experiments in this paper focused on applications of the conformance method to streamed data. It would be interesting to study how the method works on other types of vector data.

Broader Impact

As with any anomaly detection method, there might be some intrinsic ethical issues depending on the data that is used and the intended use case, particularly if it involves people. However, the authors cannot identify any ethical issues that are specific to this method.

Acknowledgments and Disclosure of Funding

This work was supported by the Defence and Security Programme at the Alan Turing Institute, funded by the UK Government. PF, TL, IPA were supported by the Alan Turing Institute under the EPSRC grant EP/N510129/1 and the EPSRC program grant EP/S026347/1 DATASIG.

The authors are grateful for the UEA & UCR time series classification repository [3], without which it would have been much more difficult to validate our approach.

References

- [1] R. J. Adler and J. E. Taylor. Gaussian inequalities. *Random Fields and Geometry*, pages 49–64, 2007.
- [2] M. Amer, M. Goldstein, and S. Abdennadher. Enhancing one-class support vector machines for unsupervised anomaly detection. In *Proceedings of the ACM SIGKDD Workshop on Outlier Detection and Description*, pages 8–15, 2013.
- [3] A. Bagnall, J. Lines, W. Vickers, and E. Keogh. The UEA & UCR time series classification repository, 2020. www.timeseriesclassification.com, accessed May 2020.
- [4] L. Beggel, B. X. Kausler, M. Schiegg, M. Pfeiffer, and B. Bischl. Time series anomaly detection based on shapelet learning. *Computational Statistics*, 34(3):945–976, 2019.
- [5] M. M. Breunig, H.-P. Kriegel, R. T. Ng, and J. Sander. LOF: Identifying density-based local outliers. In *Proceedings of the 2000 ACM SIGMOD International Conference on Management of Data*, pages 93–104, 2000.
- [6] R. Chalapathy and S. Chawla. Deep learning for anomaly detection: A survey. *arXiv preprint arXiv:1901.03407*, 2019.
- [7] V. Chandola, A. Banerjee, and V. Kumar. Anomaly detection: A survey. *ACM computing surveys (CSUR)*, 41(3):1–58, 2009.
- [8] D. Dua and C. Graff. UCI machine learning repository, 2017.
- [9] G. Flint, B. Hambly, and T. Lyons. Discretely sampled signals and the rough Hoff process. *Stochastic Processes and their Applications*, 126(9):2593–2614, 2016.
- [10] M. Hauskrecht, I. Batal, M. Valko, S. Visweswaran, G. F. Cooper, and G. Clermont. Outlier detection for patient monitoring and alerting. *Journal of Biomedical Informatics*, 46(1):47–55, 2013.
- [11] Z. He, X. Xu, and S. Deng. Discovering cluster-based local outliers. *Pattern Recognition Letters*, 24(9-10):1641–1650, 2003.
- [12] T. Hofmann, B. Schölkopf, and A. J. Smola. Kernel methods in machine learning. *The Annals of Statistics*, pages 1171–1220, 2008.
- [13] R. J. Hyndman, E. Wang, and N. Laptev. Large-scale unusual time series detection. In *Proceedings of the 2015 IEEE International Conference on Data Mining Workshop (ICDMW)*, pages 1616–1619, 2015.
- [14] A. K. Jones and R. S. Sielken. Computer system intrusion detection: A survey. Technical report, University of Virginia, 2000.
- [15] J. Kalsi, T. Lyons, and I. P. Arribas. Optimal execution with rough path signatures. *SIAM Journal on Financial Mathematics*, 11(2):470–493, 2020.
- [16] P. Kidger, P. Bonnier, I. P. Arribas, C. Salvi, and T. Lyons. Deep signature transforms. In *Proceedings of the Advances in Neural Information Processing Systems Conference*, pages 3099–3109, 2019.

- [17] F. T. Liu, K. M. Ting, and Z.-H. Zhou. Isolation forest. In *Proceedings of the 2008 Eighth IEEE International Conference on Data Mining*, pages 413–422. IEEE, 2008.
- [18] F. T. Liu, K. M. Ting, and Z.-H. Zhou. Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 6(1):1–39, 2012.
- [19] T. J. Lyons, M. Caruana, and T. Lévy. *Differential equations driven by rough paths*. Springer, 2007.
- [20] T. Mikolov, K. Chen, G. Corrado, and J. Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [21] K. Nian, H. Zhang, A. Tayal, T. Coleman, and Y. Li. Auto insurance fraud detection using unsupervised spectral ranking for anomaly. *Journal of Finance and Data Science*, 2(1):58–75, 2016.
- [22] J. F. Reizenstein and B. Graham. Algorithm 1004: The iisignature library: Efficient calculation of iterated-integral signatures and log signatures. *ACM Transactions on Mathematical Software (TOMS)*, 46(1):1–21, 2020.
- [23] Y. Wu, H. Ni, T. J. Lyons, and R. L. Hudson. Signature features with the visibility transformation. *arXiv preprint arXiv:2004.04006*, 2020.

A Plots of conformance distances for PenDigits data set

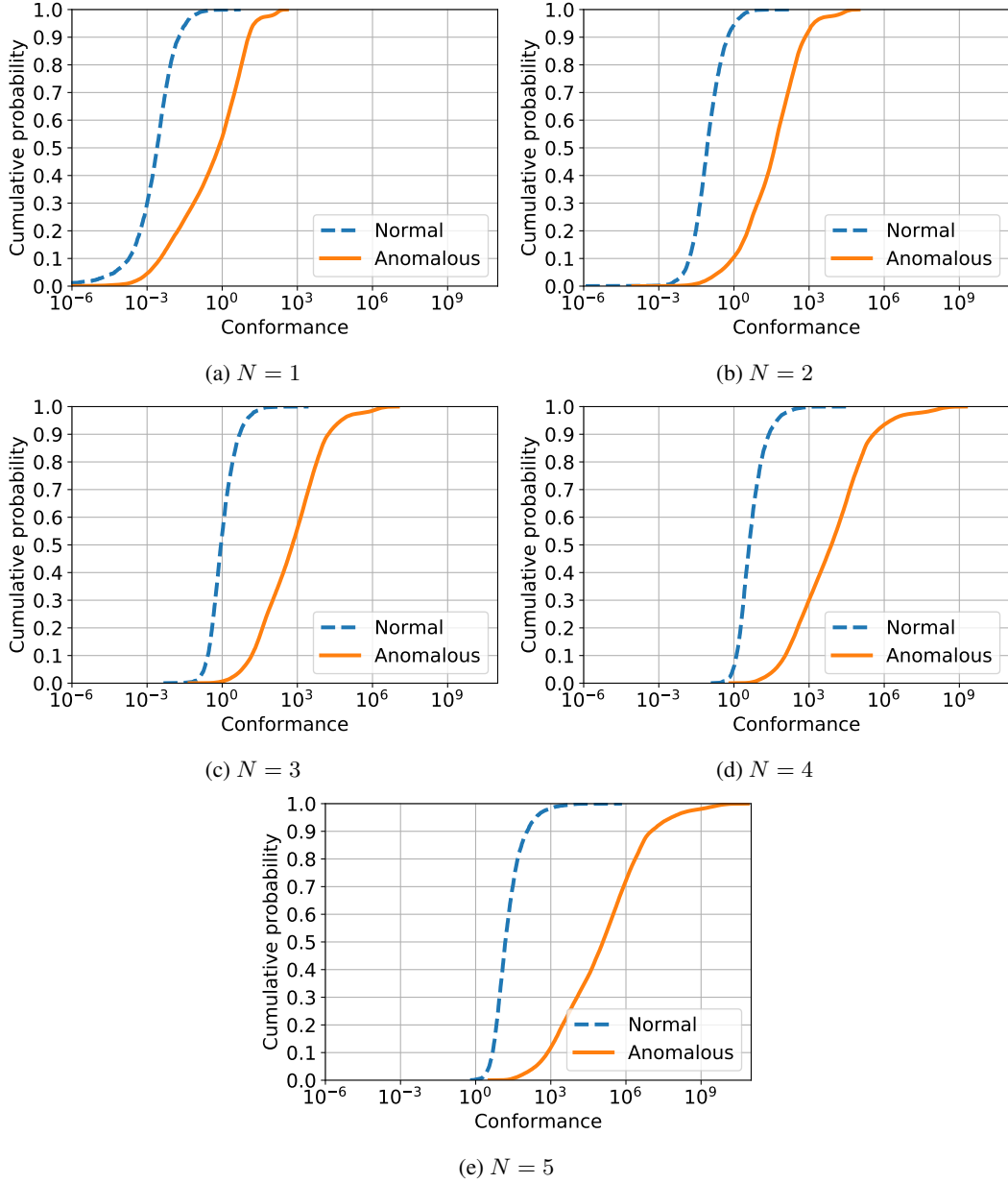


Figure 3: Empirical cumulative distributions of the conformance distance $\text{dist}(\cdot; \text{Sig}^N(\mathcal{C}))$, obtained for normal and anomalous testing data and based on computing signatures of order N .

B Table of results for univariate time series data

Data set	0.1% anomaly rate		5% anomaly rate	
	Conformance	ADSL	Conformance	ADSL
Adiac	1.00 (0.00)	0.99 (0.10)	0.99 (0.09)	0.95 (0.05)
ArrowHead	0.80 (0.07)	0.65 (0.03)	0.74 (0.06)	0.64 (0.03)
Beef	0.80 (0.22)	0.57 (0.15)	0.80 (0.22)	0.73 (0.12)
BeetleFly	0.75 (0.08)	0.90 (0.08)	0.72 (0.08)	0.84 (0.08)
BirdChicken	0.75 (0.13)	0.85 (0.15)	0.77 (0.15)	0.79 (0.09)
CBF	0.97 (0.01)	0.80 (0.04)	0.86 (0.03)	0.68 (0.03)
ChlorineConcentration	0.91 (0.01)	0.50 (0.00)	0.88 (0.01)	0.47 (0.01)
Coffee	0.80 (0.05)	0.84 (0.04)	0.78 (0.05)	0.73 (0.05)
ECG200	0.80 (0.06)	0.50 (0.03)	0.75 (0.05)	0.47 (0.04)
ECGFiveDays	0.97 (0.02)	0.94 (0.11)	0.83 (0.02)	0.86 (0.01)
FaceFour	0.78 (0.10)	0.94 (0.10)	0.78 (0.13)	0.88 (0.11)
GunPoint	0.85 (0.05)	0.75 (0.03)	0.81 (0.05)	0.68 (0.04)
Ham	0.52 (0.04)	0.50 (0.02)	0.52 (0.04)	0.50 (0.03)
Herring	0.58 (0.06)	0.52 (0.02)	0.57 (0.04)	0.49 (0.04)
Lightning2	0.73 (0.04)	0.63 (0.07)	0.75 (0.05)	0.50 (0.07)
Lightning7	0.94 (0.09)	0.73 (0.11)	0.82 (0.09)	0.68 (0.07)
Meat	0.94 (0.03)	1.00 (0.04)	0.79 (0.07)	0.87 (0.05)
MedicalImages	0.97 (0.03)	0.90 (0.03)	0.95 (0.04)	0.83 (0.05)
MoteStrain	0.89 (0.01)	0.74 (0.01)	0.86 (0.02)	0.71 (0.03)
Plane	1.00 (0.00)	1.00 (0.04)	1.00 (0.04)	1.00 (0.04)
Strawberry	0.92 (0.01)	0.77 (0.03)	0.88 (0.01)	0.67 (0.02)
Symbols	1.00 (0.01)	0.96 (0.02)	0.99 (0.01)	0.95 (0.03)
ToeSegmentation1	0.77 (0.03)	0.95 (0.01)	0.76 (0.05)	0.84 (0.03)
ToeSegmentation2	0.80 (0.06)	0.88 (0.02)	0.77 (0.06)	0.80 (0.10)
Trace	1.00 (0.00)	1.00 (0.04)	1.00 (0.05)	1.00 (0.02)
TwoLeadECG	0.92 (0.02)	0.89 (0.01)	0.82 (0.02)	0.81 (0.02)
Wafer	0.97 (0.02)	0.56 (0.02)	0.81 (0.03)	0.53 (0.01)
Wine	0.85 (0.06)	0.53 (0.02)	0.81 (0.09)	0.53 (0.02)

Table 3: Comparison of balanced accuracy from Section 4.3. Values in brackets are standard deviations with respect to testing folds.

C Properties of the variance norm for signatures

Below we give a few properties of the variance norm (1) for streamed data. Intuitively, these properties are interpreted as follows. The order of the signature $N \in \mathbb{N}$ can be seen as a measure of the resolution at which the streams are viewed. If N is small, only general features of the streams are considered. If N is increased, more and more details of the streams are considered, as they're viewed at a higher resolution.

Given a finite corpus $\mathcal{C} \subset \mathcal{S}(\mathbb{R}^d)$, any stream not belonging to the corpus is, in a way, an anomaly. In other words, viewed at a sufficiently high resolution any stream that is not in the corpus is an anomaly. The degree to which it should be considered as an anomaly should also increase with N :

Proposition C.1. *Let $\mathcal{C} \subset \mathcal{S}(\mathbb{R}^d)$ be a finite corpus. Take $w \in \mathbb{R}^{d_N}$. Then, $\|w\|_{\text{Sig}^N(\mathcal{C})}$ is non-decreasing as a function of N .*

Proof. Let $M \geq N$. We have

$$\|w\|_{\text{Sig}^N(\mathcal{C})} = \sup_{f \in \mathbb{R}^{d_N} \setminus \{0\}} \frac{\langle f, w \rangle^2}{\langle f, w \rangle^2, \mathbb{E} [\text{Sig}^{2N}(\mathbf{x})]} \leq \sup_{f \in \mathbb{R}^{d_M} \setminus \{0\}} \frac{\langle f, w \rangle^2}{\langle f, w \rangle^2, \mathbb{E} [\text{Sig}^{2M}(\mathbf{x})]} = \|w\|_{\text{Sig}^M(\mathcal{C})},$$

as the supremum is taken over a larger set. \square

Moreover, for a sufficiently high resolution, any stream of data not belonging to the corpus has infinite variance:

Proposition C.2. *Let $\mathcal{C} \subset \mathcal{S}(\mathbb{R}^d)$ be a finite corpus. Let $\mathbf{y} \in \mathcal{S}(\mathbb{R}^d)$ be a stream of data that does not belong to the corpus, $\mathbf{y} \notin \mathcal{C}$. Then, there exists N large enough such that*

$$\|\text{Sig}^n(\mathbf{y})\|_{\text{sig}^n(\mathcal{C})} = \infty \quad \forall n \geq N.$$

Proof. If $\mathbf{y} \in \mathcal{C}$, there exists N large enough such that $\text{Sig}^N(\mathbf{y})$ is independent to $\text{Sig}^N(\mathcal{C})$ [19]. Therefore, there exists $f \in \mathbb{R}^{d_N}$ such that $\langle f, \text{Sig}^N(\mathbf{x}) \rangle = 0$ for all $\mathbf{x} \in \mathcal{C}$ and $\langle f, \text{Sig}^N(\mathbf{y}) \rangle = 1$. It then follows that

$$\|\text{Sig}^n(\mathbf{y})\|_{\text{sig}^n(\mathcal{C})} = \infty \quad \forall n \geq N.$$

□