

# **Signature Isolation Forest**

Research Organization of Information and Systems The Institute of Statistical Mathematics

## Marta Campi<sup>1</sup>, Guillaume Staerman<sup>2</sup>, Gareth W. Peters<sup>3</sup>, **Tomoko Matsui**<sup>4</sup> Innia

<sup>1</sup>Institut Pasteur, Institut de l'Audition, IHU reConnect, Paris, France <sup>2</sup>INRIA, CEA, Univ. Paris-Saclay, France

<sup>3</sup>Department of Statistics & Applied Probability, University of California Santa Barbara, USA <sup>4</sup>Institute of Statistical Mathematics, Tokyo, Japan

#### 1. Context

**Functional Anomaly Detection Setting:** 

- Multivariate functional data  $\mathbf{X} \in \mathcal{F}^d$ such that  $X(t) \in \mathbb{R}^d$ ,  $\forall t \in [0, 1]$ .
- Data observed  $\mathcal{X}$ as  $\{x_i(t_1),\ldots,x_i(t_p)\}_{i=1}^n$  on a finite discretization.

#### 4. Kernel Truncated Signature 7. Parameter Analysis

The signature can be seen as a feature map<sup>3</sup> that

embeds a function or a path into the tensor algebra. It is defined as the map  $K^k : \mathcal{F}^d \times \mathcal{F}^d \to \mathbb{R}$  such that

 $K^{k}(\mathbf{X}, \mathbf{Y}) = \langle S^{k}(\mathbf{X}), S^{k}(\mathbf{Y}) \rangle.$ 

It captures non-linear relationships by embedding

paths in higher-dimensional space, preserving



AUC vs. number of split windows: isolated anomalies *(left) and persistent anomalies (right)* 

Functional Isolation Forest<sup>1</sup> (FIF) has limitations: its inner product and dictionary choices significantly impact performance. We introduce Signature Isolation Forest leveraging rough path theory<sup>2</sup>.



## 2. Contributions

Two new Functional Anomaly Detection (FAD) algorithms based on isolation forest structure and the signature approach:

- Kernel Signature Isolation Forest (K-SIF): Extension of Functional Isolation Forest (FIF) based on the kernel signature.
- Signature Isolation Forest (SIF): Isolation Forest based algorithm relying on the coordi-

sequential information and allowing comparison through geometric properties. This formulation effectively detects complex patterns in functional data.

#### 5. Kernel-SIF

**Input:** Subsample  $\{\mathbf{x}_{i_j}\}_{j=1}^m$ , dictionary  $\mathcal{D}$ , measure  $\nu$ , signature level k, split windows  $\omega$ . (a) Root node (0, 0) corresponds to  $C_{0,0} = \mathcal{F}^d$ . (b) If node (p,q) is terminal, stop; otherwise go to (c).

(c) Split node (p, q) as follows:

- 1. Choose **d** from  $\mathcal{D}$  according to  $\nu$ .
- 2. Select  $\gamma$  uniformly from

 $\min_{\mathbf{x}\in\mathcal{X}_{p,q}}\langle S^k(\mathbf{x}), S^k(\mathbf{d})\rangle,$  $\max_{\mathbf{x}\in\mathcal{X}_{p,q}}\langle S^k(\mathbf{x}), S^k(\mathbf{d})\rangle$ 

3. Form children subsets and datasets:

**Key Findings:** 

- Split Windows ( $\omega$ ):
  - Critical for isolated anomalies performance improves with more windows
  - Less important for persistent anomalies stable across window counts
  - Optimal setting  $\omega = 10$  balances performance and computational cost
- **Truncation Level (***k***)**:
  - k = 1: Captures only displacements (often insufficient)
  - k = 2: Includes path area information (good balance)
  - k > 2: Diminishing returns relative to computational cost

## 8. Real-World Performance

Dataset	SIF	K-SIF <sub>C</sub>	K-SIF <sub>B</sub>	FIF <sub>B</sub>
Chinatown	1.00	0.99	1.00	0.83
SonyRobotAI1	0.99	0.95	0.95	0.76
SonyRobotAI2	0.93	0.92	0.93	0.84
ECGFiveDays	0.93	0.92	0.90	0.93
ECG5000	0.90	0.97	0.91	0.88

nate signature.

### 3. Signature Method

Let  $\mathbf{X} \in \mathcal{F}^d$  be a r.v. of bounded variation. For any set of coordinates  $\{i_1, \ldots, i_k\} \subset \{1, \ldots, d\}^k, k \in$  $\mathbb{N}_*$ , and  $[s,t] \subset [0,1]$ , the associated **coordinate signature**<sup>4</sup> is defined by:

$$S_{(i_1,\ldots,i_k)}(\mathbf{X})_{[s,t]} = \int \cdots \int dX_{u_1}^{i_1} \ldots dX_{u_k}^{i_k}$$

$$s \le u_1 < \ldots < u_k \le t$$
(1)

Given an order of truncation  $k \in \mathbb{N}_*$ , the **truncated signature** is the vector of finite length:

 $S^k(\mathbf{X}) = (1, S_1(\mathbf{X}), \dots, S_d(\mathbf{X}),$  $S_{(1,1)}(\mathbf{X}), \dots, S_{(d,\dots,d)}(\mathbf{X}) \in \mathbb{R}^C \quad (2)$ where  $C = |\{1, ..., d\}^{\leq k}|$ 



$$C_{p+1,2q} = C_{p,q} \cap C_{\text{K-SIF}}^{L}$$
$$C_{p+1,2q+1} = C_{p,q} \cap C_{\text{K-SIF}}^{R}$$
$$\mathcal{X}_{p+1,2q} = \mathcal{X}_{p,q} \cap \mathcal{C}_{p+1,2q}$$
$$\mathcal{X}_{p+1,2q+1} = \mathcal{X}_{p,q} \cap \mathcal{C}_{p+1,2q+1}$$

(d) Apply steps (b)-(c) to nodes (p + 1, 2q) and (p+1, 2q+1)**Output:** Partition ( $C_{0,0}, C_{1,1}, \ldots$ )

## 6. Signature Isolation Forest

**Input:** Subsample  $\{\mathbf{x}_{i_j}\}_{j=1}^m$ , signature level k, split windows  $\omega$ .

(a) Root node (0, 0) corresponds to  $C_{0,0} = \mathcal{F}^d$ . (b) If node (p,q) is terminal, stop; otherwise go to (c).

(c) Split node (p, q) as follows:

- 1. Choose coordinate  $(i_1, \ldots, i_\ell)$  uniformly from  $\{(i_1, \ldots, i_\ell) \in [\![1, d]\!]^\ell; \quad 1 \le \ell \le k\}.$



*Geometric visualization of depth-2 signature terms, where*  $S^{(1,2)}$  (cyan region) and  $S^{(2,1)}$  (purple region) represent areas corresponding to coordinate signatures.

References: 1. Staerman, G., et al. (2019). Functional isolation forest. In Asian Conference on Machine Learning, pages 332–347. PMLR. 2. Lyons, T., et al. (2007). Differential equations driven by rough paths. Springer. 3. Kiraly, F. J., and Oberhauser, H. (2019). Kernels for sequentially ordered data. Journal of Machine Learning Research, 20. 4. Fermanian, A. (2021). Embedding and learning with signatures. Computational

Statistics and Data Analysis, 157:107148.

2. Select  $\gamma$  uniformly from

 $\left|\min_{\mathbf{x}\in\mathcal{X}_{p,q}} S_{(i_1,\ldots,i_\ell)}(\mathbf{x}), \max_{\mathbf{x}\in\mathcal{X}_{p,q}} S_{(i_1,\ldots,i_\ell)}(\mathbf{x})\right|$ 

- 3. Form children subsets and datasets:
  - $\mathcal{C}_{p+1,2q} = \mathcal{C}_{p,q} \cap C_{\mathrm{SIF}}^L$  $\mathcal{C}_{p+1,2q+1} = \mathcal{C}_{p,q} \cap C_{\mathrm{SIF}}^R$  $\mathcal{X}_{p+1,2q} = \mathcal{X}_{p,q} \cap \mathcal{C}_{p+1,2q}$  $\mathcal{X}_{p+1,2q+1} = \mathcal{X}_{p,q} \cap \mathcal{C}_{p+1,2q+1}$

(d) Apply steps (b)-(c) to nodes (p + 1, 2q) and (p+1, 2q+1)**Output:** Partition ( $C_{0,0}, C_{1,1}, \dots$ )

StarLightC. 500s 13-20s 202s

Signature methods are 10-30× faster than deep learning

#### 9. Conclusion

- Novel Methods: Two functional anomaly detection algorithms using signatures
- Key Advantages: K-SIF offers non-linear projections; SIF provides dictionary-free detection; both excel with sequential patterns
- **Results**: SIF achieves top performance on 50% of datasets; both methods show 5-10× speed improvement over deep learning
- **Impact**: Robust solution for functional anomaly detection with temporal ordering and non-linear patterns