

# Online Nonparametric Anomaly Detection in High-Dimensional Datasets

**Yasin Yilmaz**

Electrical Engineering  
University of South Florida

2/13/2018

# Outline

- 1 Introduction
- 2 Background
- 3 Online Nonparametric Anomaly Detection
- 4 Numerical Results
- 5 Conclusion

# Introduction

# Anomaly Detection

- **Objective:** identify patterns that deviate from a nominal behavior
- **Applications:** cybersecurity, quality control, fraud detection, fault detection, health care, . . .

# Anomaly Detection

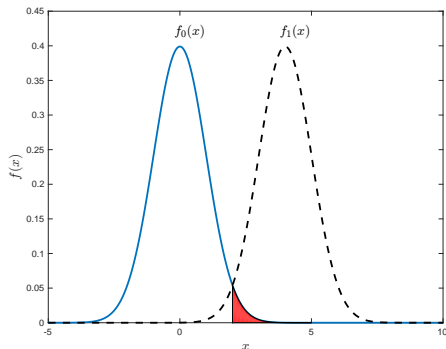
- **Objective:** identify patterns that deviate from a nominal behavior
- **Applications:** cybersecurity, quality control, fraud detection, fault detection, health care, ...

In literature typically

*statistical outlier detection*  
 =  
*anomaly detection*

However an outlier could be

- nominal tail event  
or
- real anomalous event  
(e.g., mean shift)



# Problem Formulation

Instead of *anomaly = outlier*, consider also temporal dimension

## Proposed Model

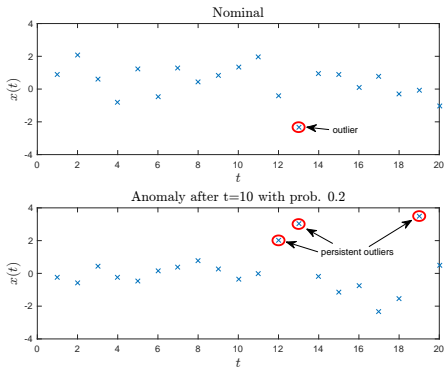
*anomaly = persistent outliers*

## Objective

**Timely** and **accurate** detection of anomalies in **high-dimensional** datasets

## Approach

*Sequential & Nonparametric* anomaly detection



## Motivating Facts: IoT Security, Smart Grid, ...

- **IoT devices:** 8.4B in 2017 and expected to hit 20B by 2020 <sup>1</sup>
- **IoT systems:** highly vulnerable – needs scalable security solutions <sup>2</sup>
- **Mirai IoT botnet:** largest recorded DDoS attack with at least 1.1 Tbps bandwidth (Oct. 2016) <sup>2</sup>
- **Persirai IoT botnet** targets at least 120,000 IP cams (May 2017) <sup>3</sup>
- **A plausible cyberattack against the US grid:** 100M people may be left without power with up to \$1 trillion of monetary loss <sup>4</sup>

---

<sup>1</sup>R. Minerva, A. Biru, and D. Rotondi, "Towards a definition of the Internet of Things (IoT)," IEEE Internet Initiative, no. 1, 2015.

<sup>2</sup>E. Bertino and N. Islam, "Botnets and Internet of Things Security," Computer, vol. 50, no. 2, pp. 76-79, Feb. 2017.

<sup>3</sup>Trend Micro, "Persirai: New Internet of Things (IoT) Botnet Targets IP Cameras", May 9, 2017, available online

<sup>4</sup>Trevor Maynard and Nick Beecroft, "Business Blackout," Lloyd's Emerging Risk Report, p. 60, May 2015.

# Motivating Facts: IoT Security, Smart Grid, ...

## Challenges:

- **Unknown anomalous distribution:** parametric methods, as well as signature-based methods (e.g., antivirus) are not feasible
- **High-dimensional problems:** even nominal distribution is difficult to know
- **Nonparametric methods** are needed
- **Timely and accurate** detection is critical



# Problem Definition

- Monitor a system online through sequential observations  $\mathcal{X}_t = \{X_1, X_2, \dots, X_t\}$  of  $d$ -dimensional independent vectors  $X_t$
- Consider an anomaly as **persistent outliers** in the observations

## Objective

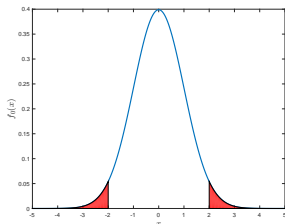
Accurately detect such anomalies in a **timely** fashion using a practical **nonparametric** approach applicable to **high-dimensional** datasets (e.g., big data problems)

# Background

# Outlier Detection

- Needs to know a statistical description of the nominal (e.g., no attack) behavior (baseline)
- Determines instances that significantly deviate from the baseline
- With  $f_0$  completely known,  $x$  is outlier if  $f_0(x) < \alpha$
- Equivalently, if  $x$  is outside the most compact set of data points under  $f_0$  (**minimum volume set**)

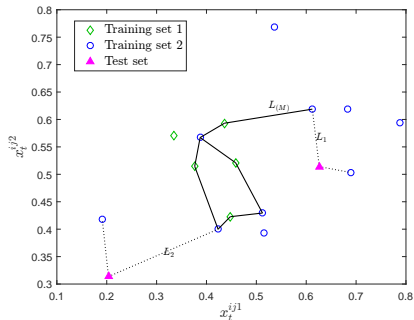
$$\Omega_\alpha = \arg \min_{\mathcal{A}} \int_{\mathcal{A}} dy \quad \text{subject to} \quad \int_{\mathcal{A}} f_0(y) dy \geq 1 - \alpha$$



- **Uniformly most powerful test** when actual distribution  $f$  is a linear mixture of  $f_0$  and the uniform distribution
- Coincides with **minimum entropy set** which minimizes the Rényi entropy while satisfying the same false alarm constraint

# Geometric Entropy Minimization (GEM)

- **High-dimensional datasets:** even if  $f_0$  is known, very computationally expensive (if not impossible) to determine  $\Omega_\alpha$
- Various methods for learning  $\Omega_\alpha$
- GEM is very **effective with high-dimensional** datasets while **asymptotically achieving** performance of  $\Omega_\alpha$  for  $\lim_{K, N \rightarrow \infty} K/N \rightarrow 1 - \alpha$



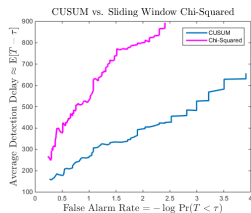
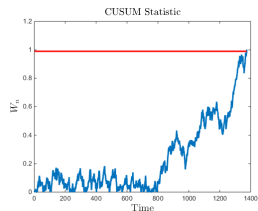
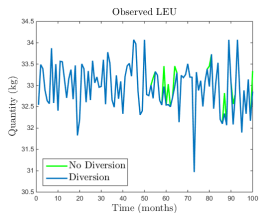
- **Training:** Randomly partitions training set into two and forms  $K$ - $k$ NN graph minimizing <sup>5</sup>

$$\bar{\mathcal{X}}_K^{N_1} = \arg \min_{\mathcal{X}_K^{N_1}} \mathcal{L}_k(\mathcal{X}_K^{N_1}, \mathcal{X}^{N_2}) = \sum_{i=1}^K \sum_{l=1}^k |e_{i(l)}|^\gamma$$

- **Test:** new point  $X_t$  outlier if  $X_t \notin \bar{\mathcal{X}}_K^{N_1+1}$

<sup>5</sup>A. O. Hero III, "Geometric entropy minimization (GEM) for anomaly detection and localization", NIPS, pp. 585-592, 2006

# Sequential Change Detection - CUSUM



$$\inf_T \sup_{\tau} \sup_{\{X_1, \dots, X_T\}} E_{\tau}[T - \tau | T \geq \tau] \quad \text{s.t.} \quad E_{\infty}[T] \geq \beta$$

$$W_t = \max \left\{ W_{t-1} + \log \frac{f_1(X_t)}{f_0(X_t)}, 0 \right\}$$

$$T = \min \{ t : W_t \geq h \}$$

# Online Nonparametric Anomaly Detection

## Online Discrepancy Test (ODIT)

- GEM lacks the **temporal aspect**

- In GEM,  $X_t$  is outlier if

$$L_t = \sum_{l=1}^k |e_{i(l)}|^\gamma > L_{(K)}$$

---

<sup>6</sup>Y. Yılmaz, “Online nonparametric anomaly detection based on geometric entropy minimization,” in IEEE International Symposium on Information Theory (ISIT), 2017, available at arXiv

## Online Discrepancy Test (ODIT)

- GEM lacks the **temporal aspect**
- In GEM,  $X_t$  is outlier if
$$L_t = \sum_{l=1}^k |e_{i(l)}|^\gamma > L_{(\kappa)}$$
- In ODIT,  $D_t = L_t - L_{(\kappa)}$  is treated as some **positive/negative evidence** for anomaly

---

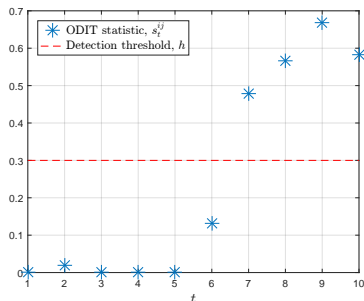
<sup>6</sup>Y. Yılmaz, “Online nonparametric anomaly detection based on geometric entropy minimization,” in IEEE International Symposium on Information Theory (ISIT), 2017, available at arXiv



# Online Discrepancy Test (ODIT)

- GEM lacks the **temporal aspect**
- In GEM,  $X_t$  is outlier if  $L_t = \sum_{l=1}^k |e_{i(l)}|^\gamma > L_{(k)}$
- In ODIT,  $D_t = L_t - L_{(k)}$  is treated as some **positive/negative evidence** for anomaly
- $D_t$  approximates  $\ell_t = \log \frac{f_1(X_t)}{f_0(X_t)}$  between  $H_1$  claiming  $X_t$  is anomalous and  $H_0$  claiming  $X_t$  is nominal
  - Assuming independence,  $\sum_{t=1}^T D_t$  gives **aggregate anomaly evidence** until time  $T$ , as  $\sum_{t=1}^T \ell_t$ , sufficient statistic for optimum detection <sup>6</sup>
  - Similar to CUSUM, optimum **sequential change detector**

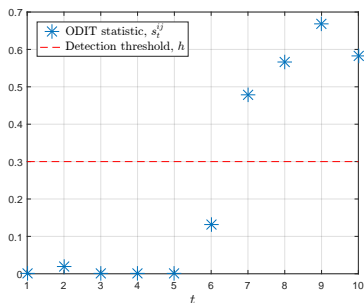
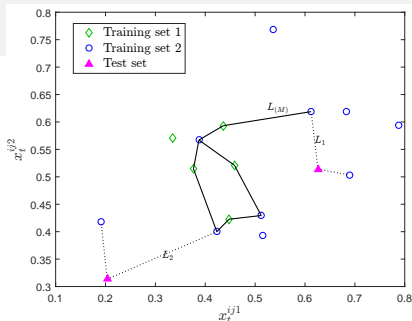
$$T_d = \min\{t : s_t \geq h\}, \quad s_t = \max\{s_{t-1} + D_t, 0\}$$



<sup>6</sup>Y. Yılmaz, "Online nonparametric anomaly detection based on geometric entropy minimization," in IEEE International Symposium on Information Theory (ISIT), 2017, available at arXiv

# ODIT Algorithm

- Initialize:  $s \leftarrow 0, t \leftarrow 1$
- Partition training set into  $\mathcal{X}^{N_1}$  and  $\mathcal{X}^{N_2}$
- Determine  $L_{(K)}$  from  $K$ -kNN graph  $\bar{\mathcal{X}}_K^{N_1}$
- While  $s < h$ 
  - Get new data  $x_t$  and compute  $D_t = L_t - L_{(K)}$
  - $s = \max\{s + D_t, 0\}$
  - $t \leftarrow t + 1$
- Declare anomaly



# Theoretical Justification - Asymptotic

## Asymptotic Optimality

- Anomaly distribution  $f_1(X_t)$  is the uniform distribution over the support of  $X_t$ :

$$H_0 : X_t \sim f_0, \forall t$$

$$H_1 : X_t \sim f_0, t < \tau, \quad \text{and} \quad X_t \sim f_{uni}, t \geq \tau$$

- as the training set grows and  $N_2 \rightarrow \infty$ ,  $D_t \rightarrow \log \frac{f_1(X_t)}{f_0(X_t)}$ ,
- and thus ODIT converges to CUSUM, which is minimax optimum.

## Sketch of the Proof

- Build an  $d$ -dimensional dynamic histogram by putting each  $k$  points in the training set in a bin.
- As the number of points increases, bin sizes decrease. Hence, the probability mass in each bin is  $k/N_2$ , where  $N_2$  is the number of all points.
- Assuming uniform distribution over bins (which is granted by the Poisson point process, however being Poisson is not required by the proof)
- probability distribution in each bin is  $(k/N_2)/(v_d r_k(X_t)^d)$  where  $v_d r_k(X_t)^d$  is the volume of the hypersphere centered at the point  $X_t$  and radius  $r_k(X_t)$
- As  $N_2 \rightarrow \infty$ , the histogram converges to  $f_0$ , thus  $(k/N_2)/(v_d r_k(X_t)^d) \rightarrow f_0(X_t)$ .
- Similarly,  $(k/N_2)/(v_d r_k(X_{(K)})^d)$  approximates a uniform distribution and as  $N_2 \rightarrow \infty$  converges to the uniform distribution over the support of  $X_t$ , where  $X_{(K)}$  is the point that corresponds to the baseline statistic  $L_{(K)}$  in the training set.

# Theoretical Justification - Nonasymptotic

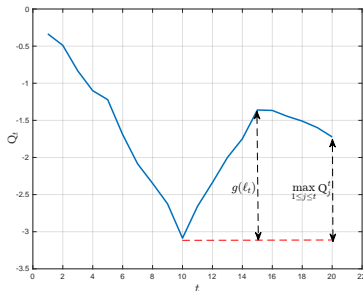
- CUSUM procedure can be expressed in terms of a general discrepancy metric, applicable to any number sequence
  - stop when discrepancy  $g(\ell_t)$ <sup>7</sup> of observations with respect to  $f_0$  is large enough

## Discrepancy and CUSUM

$$T_c = \min\{t : g(\ell_t) \geq h_c\},$$

$$\ell_t = \left[ \log \frac{f_1(X_1)}{f_0(X_1)} \dots \log \frac{f_1(X_t)}{f_0(X_t)} \right],$$

$$g(\ell_t) = \max_{1 \leq n_1 \leq n_2 \leq t} \sum_{i=n_1}^{n_2} \ell_t^i,$$



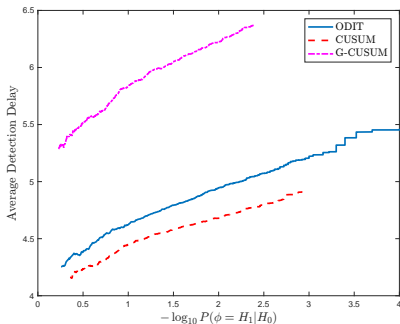
$$Q_t = \sum_{i=1}^t \ell_t$$

<sup>7</sup>B. A. Moser et al., "On stability of distance measures for event sequences induced by level-crossing sampling", IEEE Trans. Signal Process., vol. 62, no. 8, pp. 1987–1999, 2014.

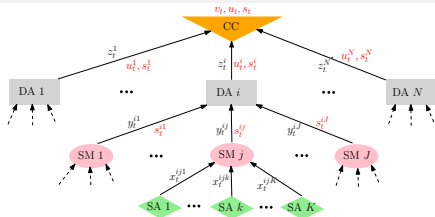
## Numerical Results

# Simulations

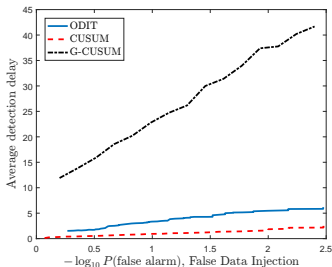
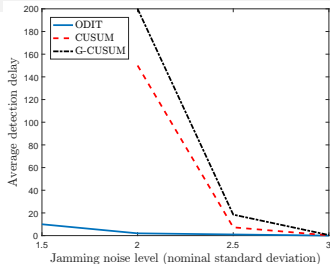
- $f_0$  is a 2D independent Gaussian with zero mean and  $\sigma = 0.1$
- $f_1 = 0.8f_0 + 0.2U[0, 1]$
- Training set 10,000 points ( $N_1 = 1000$ ,  $N_2 = 9000$ )
- $\alpha = 0.05$ ,  $k = 1$ ,  $K = \alpha N_1$
- Parametric clairvoyant CUSUM knows both  $f_0$  and  $f_1$  exactly
- Generalized CUSUM exactly knows  $f_0$ , but estimates the uniform distribution upper bound as 0.9



# Cybersecurity in Smart Grid



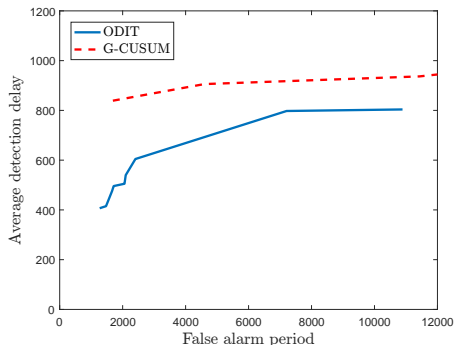
- Control center, 10 data aggregators, 1,000 smart meters, 10,000 smart appliances
- 3% of the HANs are attacked. In each attacked HAN, each smart appliance is attacked with prob. 0.5
- Baseline iid  $\sim \mathcal{N}(0.5, 0.1^2)$
- Attack data either  $\sim \mathcal{N}(0.5, (0.1\eta)^2)$ ,  $\eta > 1$  (Jamming) or  $\sim \mathcal{N}(0.5 + \Delta, 0.1^2)$ ,  $\Delta \in \mathbb{R}$  (False Data Injection)
- Even a small mismatch between the actual and assumed parameter values degrade the performance of CUSUM





# Human Activity Recognition

- Online monitoring of a dynamic system using “Heterogeneity Human Activity Recognition Dataset”<sup>8</sup> obtained from the UCI Machine Learning Repository
- Smartwatch accelerometer data: 3.5M data points with 5 numeric features
- 6 activities: biking, sitting, standing, walking, stair up, and stair down
- Focusing on activity transitions we tested online detection performance
- G-CUSUM fits multivariate Gaussian models to baseline and anomalous dist.
- Re-train after detecting a change in the activity ( $N_1 = 10$ ,  $N_2 = 20$ )



<sup>8</sup>A. Stisen et al., “Smart devices are different: Assessing and mitigating mobile sensing heterogeneities for activity recognition,” *SenSys*, 2015.

Conclusion

# Conclusions

- With the **proliferation of IoT devices**, and the **ease of triggering DoS attacks** even from unsophisticated malicious parties, there is an increasing need for developing scalable and effective solutions.
- A novel anomaly detection framework
  - **Scalable**: applicable to high-dimensional datasets (big data problems)
  - **Nonparametric**: agnostic to data-type and protocol
  - **Online** system monitoring
  - **Asymptotically optimum** for testing against uniformly distributed anomalies
- Outperforms sequential change detector CUSUM that estimates parameters from data
- Outperforms even clairvoyant CUSUM in case of a small to moderate variance increase (e.g., Jamming attack)

# Questions?

Thank you!