



Online Anomaly Detection for Optical Networks

Tarem Ahmed and Mark Coates

McGill University

tahmed@mail.mcgill.ca, coates@tsp.ece.mcgill.ca

(2)



Introduction

- High-speed optical backbones are constantly hit by wide variety of types/classes of network anomalies, from DOS attacks and viruses to large data transfers and equipment failures.
- Need online and instantaneous, anomaly detection method.
- We propose an online algorithm based on kernel version of Recursive Least-Squares (RLS) algorithm.
- Different anomalies affect network in different ways; it is not always possible to know a priori, how a potential anomaly would exhibit itself.
- Our algorithm learns characteristics of normal traffic behavior, then raises alarm immediately upon encountering a deviation from the norm.
- We test on data from the Abilene backbone network, and compare with the offline, block-based algorithm based on Principal Component Analysis (PCA) [1].

The Kernel Approach

· Data sequence is of form:

 $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), ...\}, \mathbf{x}_i \in X, y_i \in R$

where $\{\mathbf{x}_i\}_{i=1}^{t-1}$ are training vectors, and y_t is a function of $\{\mathbf{x}_i\}_{i=1}^{t-1}$.

Represent predictor as

 $\hat{y} = \sum_{i=1}^{t} \alpha_{i} \cdot \left\langle \varphi\left(\mathbf{x}_{i}\right), \varphi\left(\mathbf{x}\right) \right\rangle = \sum_{i=1}^{t} \alpha_{i} \cdot \text{kernel}\left(\mathbf{x}_{i}, \mathbf{x}\right)$

where ${\it \varphi}$ represents the mapping from input space to feature space, and α is the weight vector that minimizes the L2 error.

- Direct knowledge of *P* is not required. Only inner product of *mapped* feature vectors is required, which is provided by the chosen kernel function.
- <u>Problem</u>: once new data x_t arrives, dimension of

 $\left\{ \varphi\left(\mathbf{x}_{i}\right) \cdots \varphi\left(\mathbf{x}_{i}\right) \right\}$ may keep increasing over time, unless

- $\varphi(\mathbf{x}_{t})$ is linearly dependent on $\{\varphi(\mathbf{x})_{i}\}_{i=1}^{t-1}$.
- <u>Solution</u>: define notion of approximate linear independence as in [2]:

 $\varphi(\mathbf{x}_{t})$ is said to be *approximately* linear dependent on

(1)

 $\left\{ \varphi\left(\mathbf{x}_{i}\right) \cdots \varphi\left(\mathbf{x}_{i}\right) \right\}$, if the following condition holds:

$$\min_{a} \left\| \sum_{i=1}^{t-1} a_i \cdot \phi(\mathbf{x}_i) - \phi(\mathbf{x}_t) \right\|^2 \leq \nu$$

for some optimum sparsification coefficient vector a

Application to optical networks

- Data is normalized timeseries of number of packets or bytes, in flows between core routers in Abilene network.
- **x**_t defined to be vector giving number of packets or bytes in each flow at time *t*.
- *y_t* defined to be the total number of packets or bytes in the network at time *t*.

The Detection Algorithm

•Initialize at t = 1, by entering \mathbf{x}_1 into dictionary.

•<u>Iterate</u> for t = 2,3,...

Step 1: New dat arrives. Evaluate the approximate linear dependence of $\varphi(\mathbf{x}_{t})$ on the dictionary at time t:

$$\delta_t = \min_a \left| \sum_{j=1}^{m_{t-1}} a_j \phi(\tilde{\mathbf{x}}_j) - \phi(\mathbf{x}_t) \right|$$

Step 2: Compare δ_t with thresholds v_1 and v_2 , where $v_1 < v_2$: - If $\delta_t > v_2$, new input vector is **very** far away from space spanned by dictionary. Conclude that this is an anomaly, raise **red alarm**. No change to dictionary. - If $v_1 < \delta_t < v_2$, new input vector not sufficiently explained by dictionary. Add \mathbf{x}_t to dictionary, raise orange alarm. - If $\delta_t < v_1$, new input vector falls into normal subspace. No alarm. No change to dictionary.



Fig. 1: Outline of online anomaly detection algorithm.

Key idea: After sufficient training, dictionary should span normal subspace.

Results

- Example: Data from 11 core Abilene routers, for 1152 5-min intervals from 0000hrs Sep. 1, 2005 to 2359hrs Sep. 4, 2005.
- Using linear kernel: kernel $(\mathbf{x}_1, \mathbf{x}_2) = \langle \mathbf{x}_1, \mathbf{x}_2 \rangle$.



Fig. 2: (a) Variation in total number of packets in network; (b) growth in *D* for various values of v_1 , with $v_2 = 6v_1$.



Fig. 3: Comparing (a) δ_i in proposed algorithm with $v_1 = 0.01$ and $v_2 = 6v_1$, with (b) energy in residual subspace using block-based PCA from [1]. Spikes represent anomalies.



<u>Fig. 4</u>: Example anomaly at t = 538. Not easily seen in (a) timeseries of packets, but (b) obvious by observing distance of **x**, from each dictionary member.



Fig. 5: Distance of \mathbf{x}_t from (a) a normal dictionary member, and (b) an anomalous input vector that was admitted to dictionary.

Conclusions and Future Work

- Algorithm is recursive, there is no need to relearn from scratch when new data arrive.
- Dictionary size (*m*) is seen to stabilize after about 500 timesteps.
- Storage requirement and complexity bounded by O(*m*²), i.e. independent of time.
- Performance comparable to accepted offline, blockbased PCA method in [1].
- Work in progress includes controlling dictionary by enabling dropping of obsolete or anomalous elements, and confirming anomaly in case of orange alarm if relevant x_t is distant from subsequent input vectors.
- Future work involves letting the data determine the thresholds ν₁ and ν₂.

Acknowledgements

R. Summerhill and M. Fullmer at Abilene provided access to raw data.

References

2004

 A. Lakhina, M. Crovella and C. Diot, Diagnosing Network-Wide Traffic Anomalies, ACM SIGCOMM, Portland, OR, August 2004.
Y. Engel, S. Mannor and R. Meir, The Kernel Recursive Least Squares Algorithm, IEEE Trans. on Sig. Proc., 52(8), pp.2275--2285,