

Anomaly Detection Using AutoEncoders: the Advanced Persistent Threats Case

Ngoc Nhu Hoang

Advised by: Prof. Talal Rahwan. Department of Computer Science, New York University Abu Dhabi

Introduction

Advanced persistent threats (APTs) are stealthy and covert cyber-attacks designed to penetrate specific computer networks through exfiltration, corruption of data, or damaging critical systems. One method for identifying APTs is to detect anomalous OS-level activities, which are not easily detectable using traditional security measures as APTs stay dormant over a long time, mimic normal operations, and constitute only a minuscule part of all operations.

This paper presents and evaluates the **AutoEncoder** and **Adversarial AutoEncoder** models trained to reconstruct normal data faithfully, highlighting anomalous data whose reconstruction errors are high. The two models show good performances on par with or surpassing classical anomaly detection algorithms. The ability to isolate suspicious activities in OS-level data is of significant importance to detecting APTs early and mitigating impacts on the system.

Methodology

The following methods use deep neural networks trained on normal data to learn some meaningful compressed representation of regularities that can be used to reconstruct the data as faithfully as possible and further distinguish anomalies.

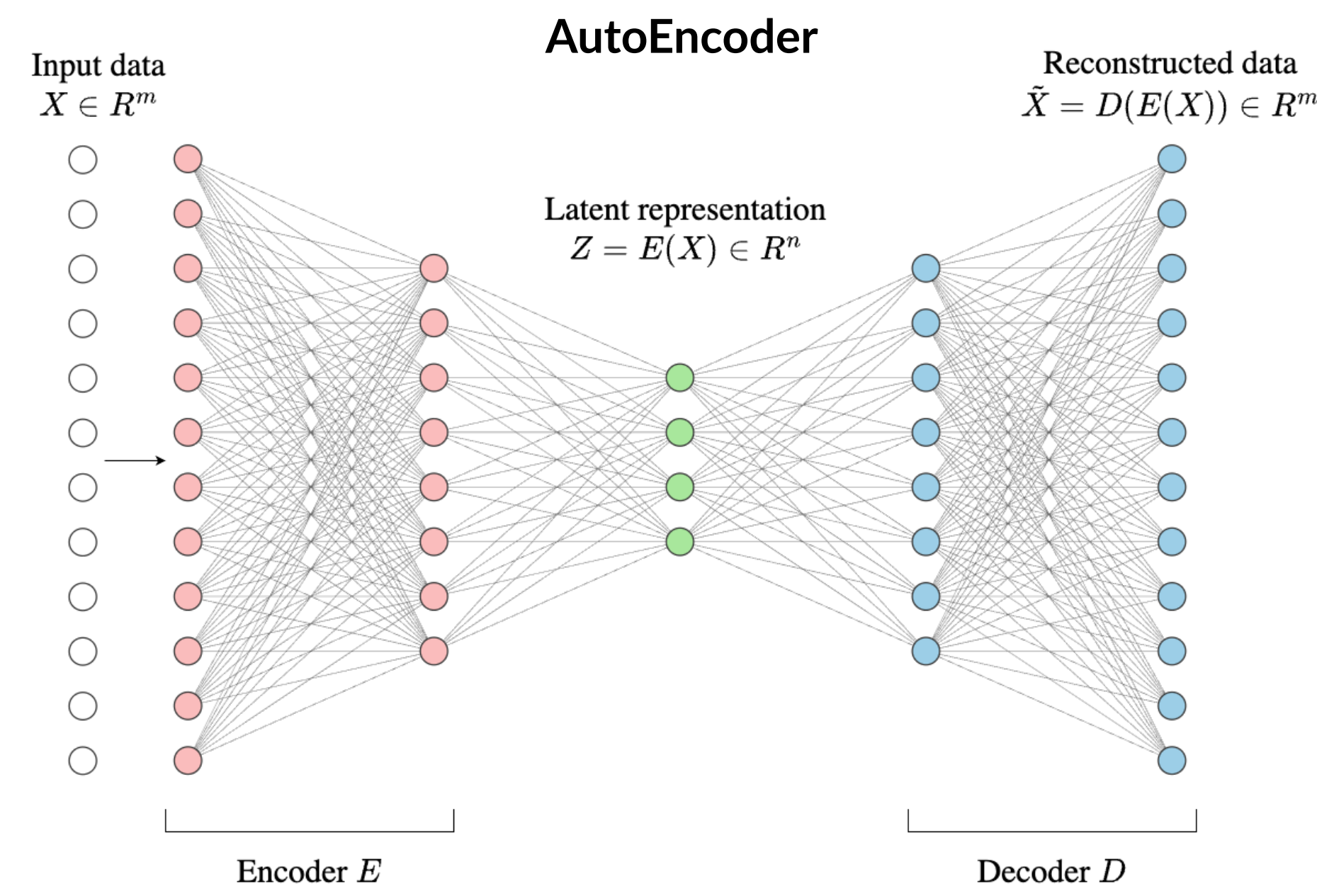


Figure 1. General architecture of the AutoEncoder model.

The AutoEncoder learns to capture the most important compressed features of the normal data that can be used to reconstruct the data as accurately as possible and will fail to reconstruct data with abnormal patterns.

Loss function: $L_{AE} = |X - \hat{X}| = |X - D(E(X))|$

Anomaly score: $A(x) = |x - \hat{x}| = |x - D(E(x))|$

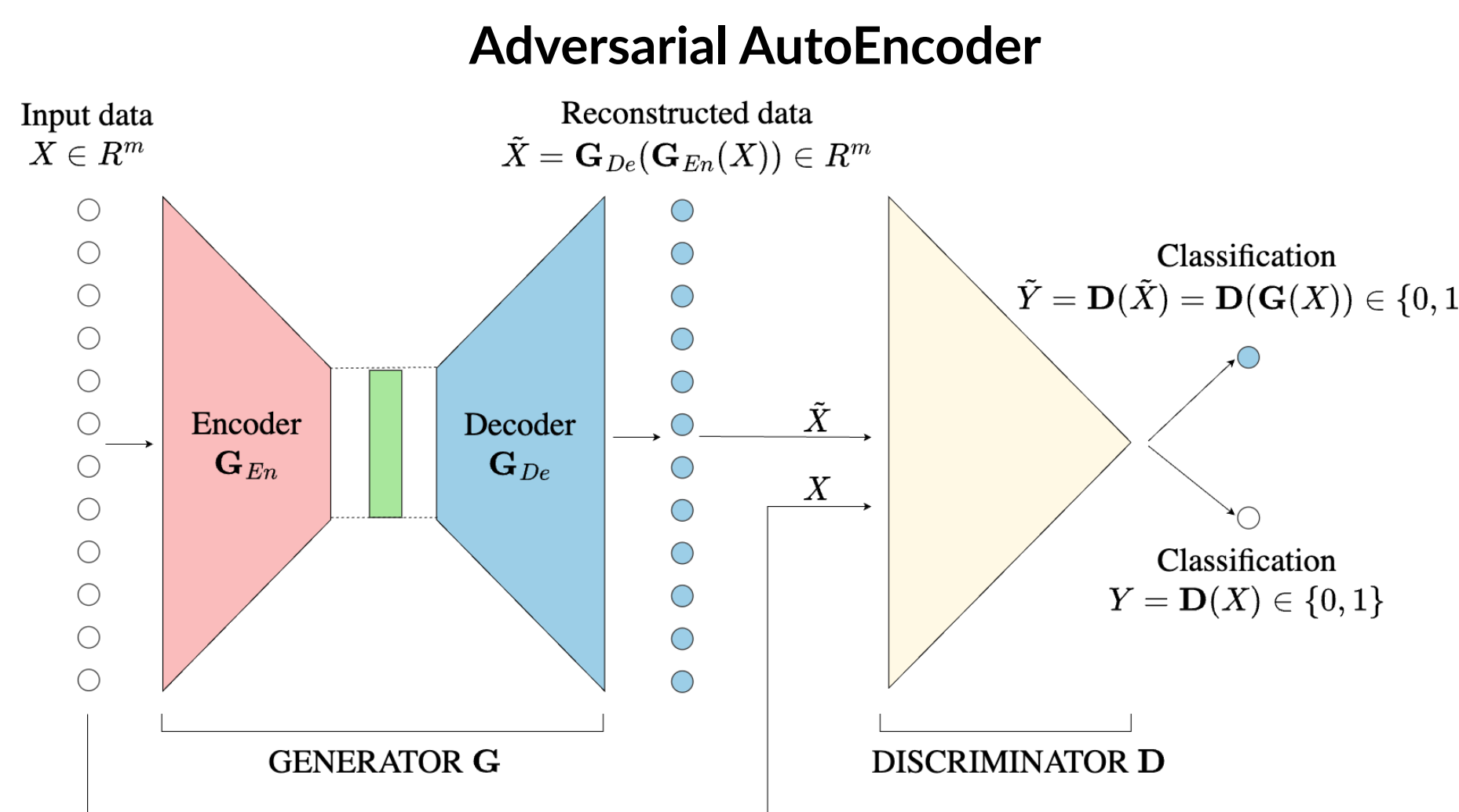


Figure 2. General architecture of the Adversarial AutoEncoder model.

Inspired by the Generative Adversarial Network [3] architecture, a deep neural network acts as the discriminator attempting to separate original data from data reconstructed by the AutoEncoder (the generator). The AutoEncoder is further motivated to reconstruct data that resemble the original data as closely as possible to fool the discriminator.

Loss functions:

- Discriminator: $L_D = |\vec{1} - \mathbf{D}(X)| + |\vec{0} - \mathbf{D}(\mathbf{G}(X))|$
- Generator: $L_G = |X - \mathbf{G}(X)| - \lambda \times L_D$

Anomaly score: $A(x) = |x - \hat{x}| = |x - \mathbf{G}(x)| = |x - \mathbf{G}_{De}(\mathbf{G}_{En}(x))|$

Data

Data from DARPA's Transparent Computing program [1] includes system operations recorded while APT-style attacks are carried out on 4 OS: **Android, Linux, BSD, Windows**. Each system has two scenarios: **Pandex** and **Bovia**. The final data includes Boolean-valued datasets called contexts, each representing an aspect of process behavior. Each OS/scenario pair has 5 contexts: **Event, Exec, Parent, Netflow, All**. In most datasets, anomalous processes constitute under 1% of all processes.

Object ID	EVENT CLONE	EVENT CHECK FILE ATTRIBUTES	EVENT OTHER	EVENT MPROTECT	EVENT CLOSE	EVENT CREATE OBJECT	EVENT LSEEK	EVENT UNLINK	EVENT WAIT	EVENT MODIFY PROCESS	EVENT RECVFROM	EVENT MODIFY FILE ATTRIBUTES	EVENT WRITE	EVENT BIND	EVENT READ	EVENT RENAME	EVENT OPEN	EVENT LOAD LIBRARY	EVENT CONNECT	EVENT SENDTO	EVENT SENDMSG
285d5fed-06dc-32ae-a04a-13cc9426616b	0	1	1	0	1	1	0	1	0	0	0	1	1	0	1	1	1	1	1	0	0
1e3548c0-b030-3591-97ac-71b67bbcb305	0	1	1	0	1	1	0	0	0	0	0	1	1	0	1	0	1	0	0	0	0
b4f1724e-0ba1-316b-973f-69e5d5e3490c	0	1	1	0	1	1	0	1	0	0	0	1	1	0	1	1	1	1	0	0	0
e2a4e818-3ce2-3626-8e22-134b542d1d77	0	0	1	0	0	0	0	0	0	0	0	0	1	0	1	0	0	0	0	0	0

Figure 3. A sample of processes from one dataset.

Evaluation method

Normalized discounted cumulative gain

The models prioritize assigning higher scores to anomalous data points and producing a ranking of the processes based on how anomalous they are. The metric chosen is the **normalized discounted cumulative gain (nDCG)** [4] typically used in Information Retrieval to measure the ability to retrieve relevant entities. nDCG score ranges from 0 to 1, with 1 being the best score i.e. when all anomalies are ranked at the top.

$$DCG = \sum_{i=1}^N \frac{rel_i}{\log_2(i+1)}; nDCG = \frac{DCG}{iDCG}$$

rel_i : relevance score of the i^{th} entry. iDCG: best achievable DCG.

Baseline algorithms

Performances of the two models are evaluated against 4 unsupervised algorithms for anomaly detection previously evaluated on the same datasets [2]: **Attribute Value Frequency, Frequent Pattern Outlier Factor, Outlier Degree, One-Class Classification by Compression**.

Results

	Process	All	Event	Exec	Parent	Netflow	
Cadets	Pandex	AE	0.8215	0.5791	0.8165	0.5358	0.1165
		AAE	0.7268	0.6394	0.6276	0.5822	0.1175
(BSD)	Bovia	AE	0.8077	0.4823	0.8168	0.7992	0.1519
		AAE	0.9079	0.4848	0.8524	0.8035	0.0760
5dir	Pandex	AE	0.5931	0.7086	0.2272	0.1970	0.6289
		AAE	0.5989	0.6676	0.2882	0.1871	0.6495
(Windows)	Bovia	AE	0.4072	0.2639	0.3256	0.4303	0.0973
		AAE	0.3952	0.2514	0.3218	0.4200	0.1298
Trace	Pandex	AE	0.6064	0.4204	0.3013	0.2397	0.3953
		AAE	0.7711	0.4916	0.4080	0.2278	0.4006
(Linux)	Bovia	AE	0.7054	0.4761	0.4871	0.2564	0.3725
		AAE	0.6494	0.4234	0.4796	0.2850	0.3633
Clearscope	Pandex	AE	0.7815	0.6708	0.4033	NA	0.6014
		AAE	0.7857	0.5484	0.5885	NA	0.6284
(Android)	Bovia	AE	0.2712	0.6357	0.5212	NA	0.3904
		AAE	0.2949	0.2631	0.3952	NA	0.1681

Table 1. nDCG scores of the two architectures on all available contexts. NA: data not available.

			All	Event	Exec	Parent	Netflow	
BSD	Pandex	AE	0.171	0.069	0.327	0.106	-0.204	-0.6 -0.4 -0.2 0.0 0.2 0.4 0.6
		AAE	0.077	0.129	0.138	0.152	-0.203	
	Bovia	AE	0.428	0.242	0.307	0.509	0.002	
		AAE	0.528	0.245	0.342	0.514	-0.074	
Windows	Pandex	AE	0.073	0.109	-0.053	-0.013	-0.081	
		AAE	0.079	0.068	0.008	-0.023	-0.060	
	Bovia	AE	0.407	0.034	0.086	0.210	0.097	
		AAE	0.395	0.021	0.082	0.200	0.130	
Linux	Pandex	AE	0.316	0.040	-0.129	-0.000	-0.085	
		AAE	0.481	0.112	-0.022	-0.012	-0.079	
	Bovia	AE	0.295	0.096	0.067	-0.164	0.372	
		AAE	0.239	0.043	0.060	-0.135	0.363	
Android	Pandex	AE	-0.048	-0.169	0.013		-0.069	
		AAE	-0.044	-0.292	0.199		-0.042	
	Bovia	AE	-0.549	0.276	0.131		-0.010	
		AAE	-0.525	-0.097	0.005		-0.232	

Figure 4. Differences in nDCG scores achieved by the models and the best-performing baseline algorithms.

Figure 4 compares the performances of the two AutoEncoder-based methods against the baseline algorithms. Each cell visualizes:

$$\Delta = nDCG_{AE/AAE} - \max(nDCG_{\{AVF, OC3, OD, FPOF\}})$$

The red cells mostly concentrate in the Android system and the **Netflow** context. These datasets either contain few processes (Android) or processes with very high dimensions (**Netflow**). In 26 out of the 38 datasets (68.4%), the best nDCG score is achieved by either the AutoEncoder or the Adversarial AutoEncoder model.

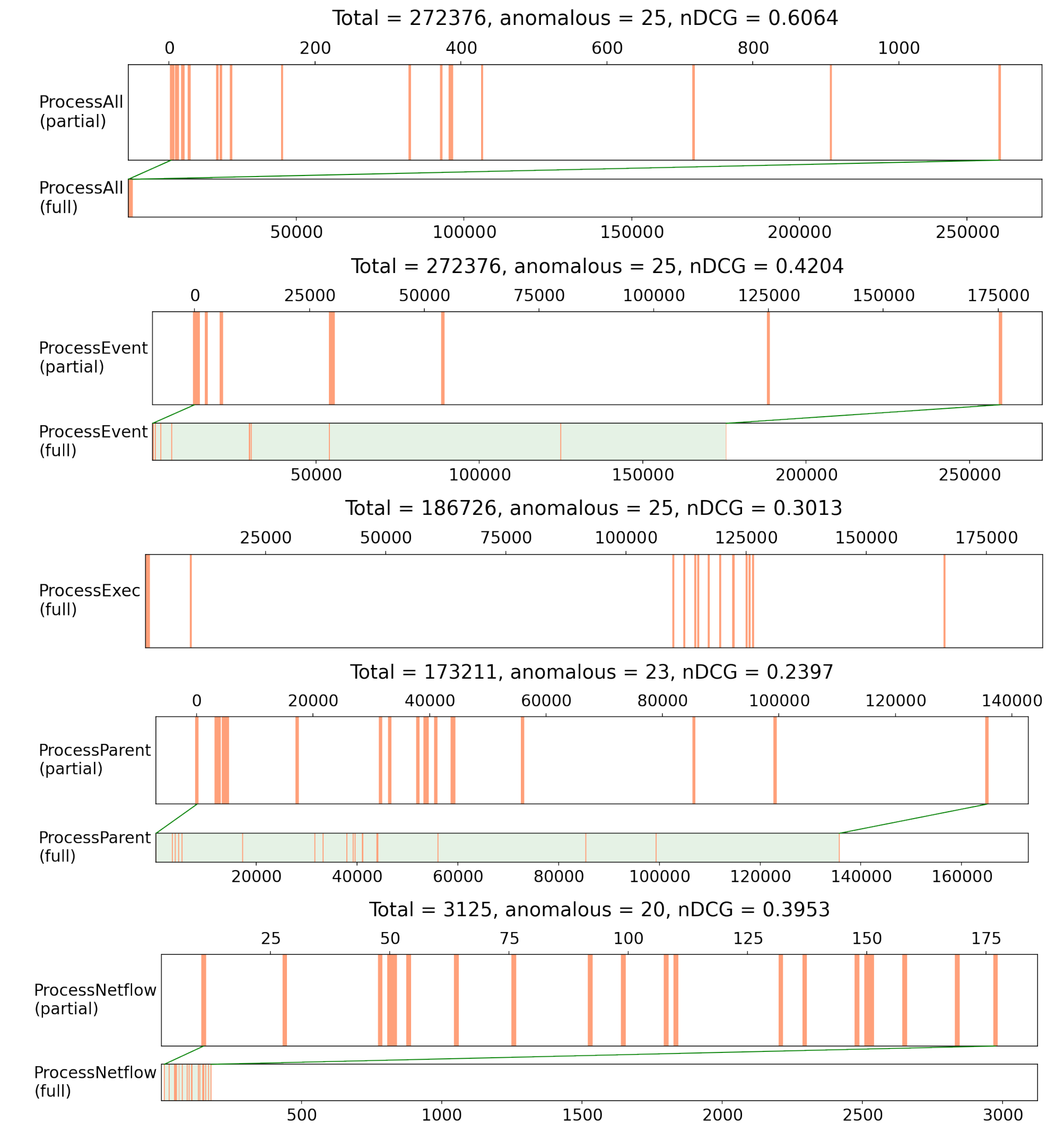


Figure 5. AutoEncoder rankings of anomalies on five contexts of the Linux system, Pandex scenario.

Conclusion and future works

This project presents two AutoEncoder-based methods to detect APTs using OS-level data. The models are trained exclusively on normal to assign anomaly scores to the processes. The two models show good performances on par with or surpassing classical anomaly detection algorithms, especially in cases where the baseline algorithms fail to finish within a reasonable time. The experiments and results from this paper can contribute to the development of efficient methods for the early detection of APTs in provenance data, potentially mitigating the damages of such attacks on the system.

In the future, different types of neural networks can be used and the layers should accommodate the varying dimensions of data. The Adversarial AutoEncoder's discriminator can be modified to better balance the workloads. Further research is needed to establish a reliable threshold for abnormality to further isolate anomalous processes.

Acknowledgement

The author thanks Professor Talal Rahwan and Sidahmed Benabderrahmane for their guidance and feedback on this project.

References

- Transparent computing.
- Ghita Berrada, Sidahmed Benabderrahmane, James Cheney, William Maxwell, Himan Mookherjee, Alec Theriault, and Ryan Wright. A baseline for unsupervised advanced persistent threat detection in system-level provenance. CoRR, abs/1906.06940, 2019.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 27. Curran Associates, Inc., 2014.
- Kalervo Järvelin and Jaana Kekäläinen. Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems*, 20(4):422–446, October 2002.