

NOTICE: this is the author's version of a work that was accepted for publication in *Expert Systems with Applications*. Changes resulting from the publishing process, such as peer review, editing, corrections, structural formatting, and other quality control mechanisms may not be reflected in this document. Changes may have been made to this work since it was submitted for publication. A definitive version was subsequently published in *Expert Systems with Applications*, [Volume 191, April 2022, 116263] DOI: 10.1016/j.eswa.2021.116263

© <2021>. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>

AutoLog: Anomaly Detection by Deep Autoencoding of System Logs

Marta Catillo^{a,*}, Antonio Pecchia^a, Umberto Villano^a

^a*Università degli Studi del Sannio, Benevento, Italy*

Abstract

The use of system logs for detecting and troubleshooting anomalies of production systems has been known since the early days of computers. In spite of the advances in the area, the analysis of log files emitted by real-life systems poses many peculiar challenges. Up-to-date tools, such as log management and Security Information and Event Management (SIEM) products, capitalize on standard data formats, logging protocols and dictionaries of threat signatures, which hardly fit to logs of industrial and proprietary systems.

This paper addresses the analysis of logs emitted by computer systems with a focus on anomaly detection. The proposed approach, named AutoLog, consists in sampling the logs at regular intervals and to compute numeric scores. Scores collected under normative operations are used to train a semi-supervised deep autoencoder, which serves as a baseline to classify future scores. The approach is not constrained by the structure of underlying logs and does not need for anomalies at training time. The results obtained in detecting anomalies of two industrial systems and the public BG/L and Hadoop datasets widely used as benchmarks, indicate that the recall of AutoLog ranges between 0.96 and 0.99, while the precision is within 0.93 and 0.98. A comparative study with isolation forest, one-class SVM, decision tree, vanilla autoencoder and variational autoencoder is conducted to demonstrate the validity of the proposal.

Keywords: system logs, deep learning, autoencoder, anomaly detection, cybersecurity

*Corresponding author (phone: +39-0824-305805)

Email addresses: marta.catillo@unisannio.it (Marta Catillo), antonio.pecchia@unisannio.it (Antonio Pecchia), villano@unisannio.it (Umberto Villano)

1. Introduction

System logs are sequences of time-stamped text lines that report on informational and abnormal events occurring during the execution of a given system. Almost any computer system appends lines to one or more special files –called *system logs*, *log files* or just *logs*, for short– at run-time. Since the early days of computers, to “grep” log files to search interesting keywords (e.g., *error*, *denied* or *unavailable*) is a well-established mean to monitor and assess dependability of computer systems and to gain direct insight into failures, anomalies and misuse (Oliner et al., 2012; Cinque et al., 2016). Nowadays, many log-related tasks have converged into up-to-date log management tools and Security Information and Event Management (SIEM) products (Miller et al., 2010; Bhatt et al., 2014), such as LogRhythm¹, Splunk² and Logstash³, which allow to collect and normalize diverse data sources and logs. More importantly, these tools and products implement real-time monitoring and alerting capabilities, which are critical for practitioners and administrators to develop situational awareness from large volumes of run-time logs. In spite of the “technical” advancements in log management, substantial cognitive work by human experts and system administrators is required to traverse the logs in order to pinpoint and to correlate relevant lines for forensics and troubleshooting.

The effectiveness of log management and SIEM installations is intertwined with the completeness of the specifications of *interesting events*, i.e., events that should be detected and followed up for further inspection. Specifications may consist in *hard-coded* catalogues of keywords, regular expressions and rules, which are used to scan system logs and to alert analysts and administrators upon matches. Figure 1 provides a concrete example of rule from OSSEC⁴, i.e., an open source log monitoring solution. The rule –coded in eXtensible Markup Language (XML)– is intended to raise an alert whenever the phrase “Client sent malformed Host header” (within the `match` tag) is seen in the logs of the Apache web server; the phrase is deemed a potential

¹<http://logrhythm.com/products/siem>

²<http://www.splunk.com/>

³<https://www.elastic.co/products/logstash>

⁴<https://ossec.github.io/docs/manual/non-technical-overview.html>

```

1 <rule id='30107' level='6'>
2   <if_sid>30101</if_sid>
3   <match>Client sent malformed Host header</match>
4   <description>Code Red attack.</description>
5   <info type="link">
6     http://www.cert.org/advisories/CA-2001-19.html
7   </info>
8   <info type="text">
9     CERT: Advisory CA-2001-19 "Code Red" Worm Exploiting
10    Buffer Overflow In IIS Indexing Service DLL
11  </info>
12  <group>automatic_attack,</group>
13 </rule>

```

Figure 1: Example of rule from OSSEC (*apache_rules.xml* file).

indicator of a *Code Red attack*. Current products rely on internal representation formats (e.g., MDI Fabric and Common Information Model used by LogRhythm and Splunk, respectively) and a variety of default *log adapters* along with comprehensive catalogs of rules for many standard protocols and commodity applications that can be encountered in production environments (e.g., ftp, ssh, telnet, squid, nginx). For other types of proprietary or industrial log files, they require users to perform ad-hoc data preparations. Moreover, default rules must be typically specialized and updated by system administrators and practitioners: setting up well-crafted and effective rules entails substantial **threat knowledge** by domain experts.

This paper addresses the analysis of logs emitted by computer systems, with a focus on the detection of anomalies, such as failures and misuse. We consider the logs of four real-life systems, which range from an **industrial system** in the transportation domain to a **microservices-based installation** implementing a standard multimedia architecture adopted by large *telcos*, up to publicly-available system logs from a Blue Gene/L (BG/L) **supercomputer** and a Hadoop **cluster**. Overall, the systems consist of various nodes and applications interacting through a network, and strongly rely on logs to massively record execution events, traces and dumps of variables. Our work stems from different motivations. There exist many classes of systems where it is hard to fit the concept of pre-established rules, such as the one shown in Figure 1. Proprietary and vendor-dependent logs, including those generated by the industrial and microservices system addressed by our study, lack standardized practices (Cinque et al., 2018). Moreover, differently from many conventional protocols and standard data sources (e.g., netflows, web

56 servers or intrusion detectors), there is not yet a mature threat model or
a default catalogue of rules for monitoring the logs of transportation and
telco systems. These problems bear the risk of underutilizing logs emitted
by real-life systems: in fact, although logs are a goldmine of information,
60 their analysis is still a great challenge. Given the volume of logs, manu-
ally inspection is hard and impractical. Anomaly detection techniques can
be conveniently used on the top of system logs to overcome the limitations
posed by the need for pre-established rules.

64 We propose **deep autoencoding of system logs** (AutoLog) to mitigate
the issues above. AutoLog consists in sampling the logs of the system under
assessment at regular intervals, and to compute numeric scores by means of a
term-weighting technique without making any specific assumption on the
68 format of underlying logs. Scores collected solely under *normative operations*
are used to train a **semi-supervised deep autoencoder**, which serves as a
baseline model to classify future scores into normative or anomaly classes. It
should be noted that semi-supervised learning does not need to know anoma-
72 lies at training time: in this respect, detection is pursued by pinpointing the
scores that deviate –through the notion of **reconstruction error** of the
autoencoder– from the baseline model. This leads to a more general solution
to log-based anomaly detection. AutoLog is suited to distributed systems
76 because it is inherently conceived to aggregate scores computed with the log
lines gathered from different sources (e.g., log files) at a given time. It is
worth noting that current distributed systems may create distinct log files
per application or node. In this respect, clock synchronization across the
80 nodes of the system is a requisite to time-stamp and to process the logs
correctly. It is well-accepted that a synchronization protocol, such as the
Network Time Protocol (NTP), provides accuracies generally in the range of
0.1 ms with fast local area networks (LANs) and computers and up to a few
84 tens of milliseconds in the intercontinental Internet⁵: this delay is negligible
when compared to the sampling period of AutoLog, which is in the range
ten to hundreds seconds. AutoLog requires no a-priori catalogues or rules of
interesting symptoms or patterns. Overall, the approach makes it possible
88 to cope with the lack of threat models for handling proprietary logs and to
complement insights from hard-coded detection rules.

AutoLog is evaluated by means of direct experiments with the logs of

⁵<https://www.eecis.udel.edu/~mills/ntp.html>

above-mentioned systems, where we compute and analyze numeric scores reflecting normative operations and anomalies. Experiments are based on a mixture of simulated and spontaneous anomalies. In the industrial and microservices systems we reproduce bruteforce authentication attempts, tampering with data structures and system misuse –inspired by accepted taxonomies in the area– and collect the logs; on the other hand, the BG/L log accounts for hardware and software errors observed over 215 days of operations that were neither induced nor simulated. The Hadoop dataset accounts for different types of service failures in the production environment. Results indicate that the recall of AutoLog at detecting anomalies ranges between 0.96 (industrial system) and 0.99 (microservices system); precision ranges between 0.93 (BG/L) and 0.98 (microservices system). As for BG/L, whose reference dataset is a widely-used benchmark in log analysis and anomaly detection, AutoLog achieves recall, precision and F1 score of 0.98, 0.93 and 0.95: these figures are strongly competitive with other existing methods assessed on the same BG/L dataset. The assessment of AutoLog is complemented by extensive discussion on the use of Principal Component Analysis (PCA) and clustering in our domain; moreover, we compare AutoLog with a wide set of techniques including isolation forest, one-class Support Vector Machine (SVM), decision trees, vanilla autoencoder and variational autoencoder.

The paper is organized as follows. Section 2 presents related work in the area. Section 3 describes the reference systems, our approach to compute scores from logs and available datasets. Section 4 addresses deep autoencoding for anomaly detection, design and training aspects. Section 5 reports the assessment of AutoLog. Section 6 proposes a comparative study of AutoLog with respect to other techniques. Section 7 discusses limitations and threats to validity of our study, and how they have been mitigated, while Section 8 concludes the work and provides future research perspectives.

2. Related Work

2.1. Mining Quantitative Metrics from Logs

Mining and analyzing quantitative **metrics** and **scores** inferred from system logs and other monitoring tools has been proven to be successful for addressing dependability and security problems in a variety of domains. A literature review on the analysis of logs for vulnerability and security is presented in (Svacina et al., 2020).

The paper (Farshchi et al., 2018) address *sporadic* Cloud operations. The technique correlates logs and Cloud metrics to detect anomalies during operations. Lines in the logs are clustered into higher-level activities through the Pearson product-moment correlation coefficient method. A regression-based technique is used to infer a correlation between lines in the log and changes in Cloud metrics. Correlations are used to formulate assertions, which aim to detect deviations of the system behavior.

Natural language processing (NLP) techniques have been used to address system logs. An anomaly detection technique leveraging NLP is proposed in (Bertero et al., 2017); the approach is intended for analyzing logs from a software system in face of different operating conditions. The analysis is performed through the Google *word2vec* algorithm, which allows mapping words into a high dimensional space. Based on the mapping, vectors of features are created to train a classifier and to provide information on the target operating condition of the system.

More recent contributions in this field attempt to obtain better representations than *word2vec* by capturing semantic information hidden in log templates. The LogAnomaly framework (Meng et al., 2019) leverages a novel word representation method based on synonyms and antonyms, *template2Vec*, to effectively represent the words in templates. The LogBERT framework for log anomaly detection (Guo et al., 2021) is instead based on Bidirectional Encoder Representations from Transformers (BERT). BERT is a Google-developed Transformer-based machine learning technique for NLP pre-training. By using the structure of BERT, LogBERT expects that the contextual embedding of each log entry can capture the information of whole log sequences.

An approach for mining console logs to detect run-time problems in large-scale systems is presented in (Xu et al., 2008). The approach extracts structured information from console logs and constructs vectors of features.

The Authors of reference (Campos et al., 2018) present a study on the use of machine learning for failure prediction. The study uses datasets of failure and non-failure data, which consist of numeric features representing the system behavior.

A log-based abnormal task detection approach for Apache Spark is presented in (Lu et al., 2017). The approach leverages a set of features extracted from the logs of Spark, e.g., execution time and data locality of each task, to detect where and when abnormalities occur.

The work (Zoppi et al., 2016) presents an anomaly detection approach

164 for Service-Oriented Architectures (SOAs), which aims to cope with SOAs’
dynamics by collecting metrics at different system layers.

The technique proposed in (Oprea et al., 2015) aims to detect early-stage
infections targeting enterprise networks. The technique leverages network
168 logs, such as Domain Name System (DNS) and web proxy logs, and com-
putes scores for domains contacted by known compromised hosts. A graph-
theoretic approach –namely *belief propagation*– is used to identify domains
that are indicative of malware infections.

172 2.2. Deep Autoencoding for Anomaly Detection

The model that we use in AutoLog to capture a normative baseline is
based on **deep autoencoding**. This was first used for anomaly detection in
(Hawkins et al., 2002), becoming progressively more and more popular in re-
176 cent years. Since then, several studies have applied autoencoders for anomaly
detection. For example, in (Sakurada and Yairi, 2014) autoencoders, denois-
ing autoencoders, PCA, and kernel PCA methods are compared according
to their performance. The use of deep learning models for anomaly detection
180 is surveyed in (Ruff et al., 2020; Pang et al., 2021).

With respect to security applications, autoencoders typically underlie
hybrid designs of anomaly detectors. In (Fahimeh and Heikkonen, 2018)
the Authors leverage a deep autoencoder to create an anomaly-based intru-
184 sion detection system (IDS), and evaluate its performance by means of the
KDD-CUP’99 dataset. In particular, their model is based on a stacked au-
toencoder, and uses in the training phase a greedy unsupervised layer-wise
training mechanism (Hinton et al., 2006).

The work (Shone et al., 2018) proposes an intrusion detection model that
is a combination of deep and shallow learning. In particular, it uses a non-
symmetric deep autoencoder (NDAE) for unsupervised feature learning, and
a classification model based on stacked NDAEs and the Random Forest al-
192 gorithm. The performance of this solution is evaluated on the KDD-CUP’99
and NSL-KDD datasets. In (Catillo et al., 2020) the Authors describe a semi-
supervised learning technique that provides two different training phases.
They test the approach on CICIDS2017 dataset by exploiting the double
196 loop learning concept with the aim of reducing the number of false positives.

Authors in (Nguyen et al., 2019) propose a framework for detecting and
explaining anomalies in network traffic. They leverage a variational autoen-
coder in order to detect anomalies. They demonstrate the validity of the

200 proposal on the recent University of Granada (UGR) dataset (Maci-Fernandez
et al., 2018). The method is effective for detecting a variety of attacks.

An autoencoder-based approach is proposed by the Authors of (Aygun
and Yavuz, 2017). In particular, they describe two deep learning-based
204 anomaly detection models using an autoencoder and a denoising autoencoder
to detect zero-day attacks.

In the context of industrial anomaly detection, it is worth mentioning the
Autoencoder-based Payload Anomaly Detection (APAD) method (Kim et al.,
208 2020). Its Authors propose two payload anomaly detection approaches by
leveraging an autoencoder: they are required at operative level and product
process management level, respectively.

In (Liu et al., 2021), instead, a semi-supervised anomaly detection method
212 is proposed, based on the encoder-decoder-encoder paradigm. The Authors
show the effectiveness of their proposal by means of an extensive experimen-
tation conducted on their Aluminum Profile Surface Defect (APSD) dataset.
Deep SAD (Ruff et al., 2020) is a deep methodology for semi-supervised
216 anomaly detection. Its Authors propose the use, in addition to the labeled
normal samples commonly exploited by semi-supervised approaches, of a
small set of labeled anomalies to obtain performance improvements.

The work (Su et al., 2019) proposes a stochastic recurrent neural network
220 approach for multivariate time series anomaly detection. It captures tempo-
ral dependence between multivariate observations and applies a variational
algorithm for representation learning. The approach aims to reconstruct
input data by the representations and use the reconstruction probabilities
224 to determine anomalies. Different from our work, experiments are done on
numeric datasets and detection achieves an overall F1 score of 0.86.

It is worth noting that autoencoders are also used in domains other than
anomaly detection. For example, paper (Zhao et al., 2021) proposes a condi-
228 tional variational autoencoder to solve a highly imbalanced classification
problem, where the training data points of the minority classes are rare. In
(Zhang et al., 2021) the Authors attempt to solve the generalized feature
selection problem. Also in this case they use an autoencoder to reduce the
232 dimensions of data while maintaining a high-quality representation as well.

2.3. Our Contribution

While deep autoencoding methods are increasingly used for anomaly and
intrusion detection in well-structured network traffic records –as for many
236 of the papers referenced in Section 2.2– we take a different perspective by

addressing text lines of system logs. The analysis of system logs to detect anomalies by “traditional” methods (i.e., not based on deep learning) is presented in (He et al., 2016b); the paper also includes the evaluation of six
240 anomaly detection methods (three supervised and three unsupervised) on publicly-available datasets. Log anomaly detection based on deep learning is instead surveyed in (Yadav et al., 2020).

Recent trends in the intersection of deep learning and log analysis put
244 forth the use of Long Short-Term Memory (LSTM) approaches. The work (Yang et al., 2019) proposes *nLSALog*, an anomaly detection framework that leverages log files as data source. The framework models the log as a natural language sequence and uses LSTM –built on the top of nominal training
248 data– to detect security anomalies. The Authors of (Yuan et al., 2020) propose Adaptive Deep Log Anomaly Detection (ADA), which aims to detect security-related anomalies in system logs, leveraging deep neural networks with LSTM and dynamic adaptive thresholds. The approach in (Du et al.,
252 2017b), named *DeepLog*, uses a deep neural network to model a system log as a natural language sequence. DeepLog learns patterns from normative executions in order to detect anomalies. Similarly, the Authors of (Zhang et al., 2016) parse streamed console logs to detect early warning signals for
256 IT system failure prediction by means of log pattern extraction.

The studies mentioned above capitalize on LSTM, which *does* allow to capture dependencies over sequences of lines in a given log file, but at the cost of *complex* and *error-prone* data preparation and clustering of similar
260 lines of the logs into common templates (or patterns). VeLog –close to our work– is an anomaly detection method based on variational autoencoders (Qian et al., 2020). VeLog needs to generate the order and number of log execution matrices: a new log sequence is labeled as an anomaly if the order
264 of log execution and number of log execution are predicted to be abnormal. As LSTM-based approaches, VeLog strongly depends on sequences of lines.

Two autoencoders along with Isolation Forest are used for unsupervised anomaly detection in logs in (Farzad and Gulliver, 2020). The same Au-
268 thors propose in (Farzad and Gulliver, 2021) the extraction of features from log messages using an autoencoder, successively exploiting an LSTM, Bidirectional Long-Short Term Memory (BLSTM) or a Gated Recurrent Unit (GRU) to classify the extracted features. The paper (Wadekar et al., 2019)
272 presents a solution utilizing a hybrid Convolutional Autoencoder-Variational Autoencoder (CAE-VAE) architecture. Keys derived from individual entries of log files are grouped in discrete event sequences, and a likelihood metric

is used as an anomaly score.

276 Differently from the majority of papers cited above, AutoLog does not
depend on the notion of sequence of lines in the logs. We address hetero-
280 geneous –and potentially distributed– logs by extracting vectors of numeric
scores, which make it possible to apply a wide category of classifiers and deep
learning models. While doing so, AutoLog embeds no application knowledge
and makes no assumptions on the format and sequences of underlying lines
in the logs. As for the computation of the metrics, differently from (Farshchi
284 et al., 2018) and (Xu et al., 2008), our proposal does not require supplemen-
tary data sources or the availability of the source code of the applications in
hand; with respect to (Bertero et al., 2017) and (Oprea et al., 2015), we do
not target a specific type of operation. Accordingly, AutoLog can be applied
and ported across different systems, as we did in our study.

288 3. Systems and Log Analysis Approach

In the following, we describe the reference systems, our log analysis ap-
proach to compute quantitative scores from logs and available datasets of
scores –from normative operations and anomalous conditions– used to con-
292 duct the experiments.

3.1. Reference Systems

Our study leverages four systems, which range from a **proprietary in-**
formation system in the transportation domain by a top industry vendor⁶,
296 to a **microservices-based installation**, up to publicly-available system
logs –common benchmarks in the literature– from a Blue Gene/L (BG/L)
supercomputer and a Hadoop **cluster**. A brief description of the systems
is provided below:

- 300 • **Industrial system.** It consists of seven nodes within a local area
network (LAN). The nodes host a variety of applications that handle
transportation-related data, such as current vehicle positions and ex-
pected routes, and implement several critical functions, which include
304 *vehicle monitoring, route computation, trajectory tracing, alerting* and
human-machine interfaces; each node hosts one or more applications.

⁶The name of the vendor is not disclosed due to confidentiality reasons.

308 Noteworthy, the system is operated with *vendor-provided* client applications, which serve as workload generator. The workload generated by the clients consists of a mixture of service requests that mimic representative usage profiles of the system in production.

- 312 • **Microservices system.** The system consists of a Clearwater⁷ microservices installation, which implements the standard IP Multimedia Subsystem (IMS) architecture being adopted by large *telcos* for IP-based voice, video and messaging services. Microservices are connected through a LAN environment and each characterized by a unique IP address. An “anchor” microservice –named *bono* in Clearwater– serves as interface to external clients; other key functions include a Session Initiation Protocol (SIP) router, a database of profile data and a RESTful server that allows authentication credentials and users profiles to be retrieved. As the previous system, Clearwater is exercised with representative usage scenarios and workloads⁸, such as registering and deleting accounts, sending SIP messages and creating endpoints.
- 320 • **Supercomputing system.** We use publicly-available system log data⁹ of a BG/L supercomputer from the Lawrence Livermore National Labs (LLNL) consisting of 131072 processors, whose reference dataset (Oliner and Stearley, 2007) has become a consolidated benchmark in the area of log analysis and anomaly detection. Interested readers are referred to (Adiga et al., 2002) for an overview of the BG/L supercomputer.
- 328 • **Hadoop system.** We use publicly-available logs¹⁰ generated by a Hadoop cluster with 46 cores running distributed applications backed by the Hadoop Distributed File System (HDFS). Hadoop is a Big Data framework that allows for the distributed processing of large data sets; it is widely used and studied in the literature given its relevant use by industry. Interested readers are referred to the reference website¹¹ for any additional information.

⁷<http://www.projectclearwater.org/>

⁸<https://github.com/metaswitch/clearwater-live-test/>

⁹<https://www.usenix.org/cfdr-data>

¹⁰<https://github.com/logpai/loghub/tree/master/Hadoop>

¹¹<http://hadoop.apache.org/>

Table 1: Examples of log lines from the reference systems.

log line	system
21 00:08:36.558 [Thread-024] Time: FunctionName: Warning: No database connection found for key '01234'	industrial
[05/21/16 00:07:11.902] Message: received=0x20 datasize=256 dataid=4	industrial
May 21 00:09:39 NODE1 ntpd[11080]: synchronized to 192.168.56.101 stratum 2	industrial (syslog)
6-11-2018 14:18:32.644 UTC Error bono.cpp:1337: Route header flow identifier failed to correlate	microservices
2005-06-03-15.42.50.363779 R02-M1-N0-C:J12-U11 RAS KERNEL INFO instruction cache parity error corrected	BG/L
2015-10-17 15:37:57,902 INFO [main] org.apache.hadoop.mapreduce.v2.job- history.JobHistoryUtils: Default file system [hdfs://msra-sa-41:9000]	Hadoop

System logs collected in distributed deployments –such as those consid-
 336 ered in our study– typically result from the “intertwinement” of heteroge-
 neous lines recorded by a variety of daemons, applications and nodes over the
 time. Log lines reporting the events may be either centralized at a unique
 location, as for BG/L, or sparse across different files. For example, our study
 340 leverages 16 log files from the industrial system and 13 log files from the
 microservices system; similarly, the Hadoop dataset consists of several log
 files populated by the system during a given timeframe. Overall, the logs
 addressed by our study encompass a mixture of formats that consist of a
 344 small number of standard fields (e.g., *time-stamp*, *hostname*, *severity*) and
 –mostly– unstructured text messages.

Table 1 provides examples of real log lines from the reference systems.
 In the industrial system we leverage both logs in vendor-dependent formats
 348 and operating system logs –typically available at `/var/log/syslog` in Linux-
 based systems– according to the widely-used `syslog` protocol¹². It is worth
 noting that the first two lines in Table 1 –although from the same vendor– re-
 veal different formats, such as the time-stamp, lack of *thread id* in the second
 352 line and a different syntax for representing key-value pairs. The microservices
 system and BG/L rely on their own log formats as well. For example, BG/L
 is based on an effective reliability, availability and serviceability (RAS) log-

¹²<https://tools.ietf.org/html/rfc5424>

ging framework via a centralized DB2 database (Oliner and Stearley, 2007);
356 lines in the log are accompanied by various flags, such as originating loca-
tion (e.g., rack and node) and severity (e.g., INFO, ERROR and FATAL). On
the other hand, Hadoop uses `log4j`¹³ to handle the logs, which introduces a
360 further log format in the context of our analysis. The inherent heterogeneity
of purposes and formats of logs that can be encountered in real-life systems
is a challenge to practitioners. In this study we propose a *uniform* analysis
and detection approach than can be ported across different types of logs.

3.2. Overview and Log Analysis Approach

364 3.2.1. Overview of AutoLog

AutoLog consists in sampling system logs with period P and then follow-
ing up with two *pipelined* steps: (i) computation of numeric scores from the
log lines, which can be handled more conveniently for machine learning than
368 the original raw text shown in Table 1, and (ii) anomaly detection based on
the scores. Figure 2 shows the overall architecture of AutoLog, where the
steps mentioned above are named *scoring* and *anomaly detection*, respec-
tively. The approach relies on a database –shown in Figure 2– of normative
372 *chunks*, i.e., log lines windows of duration P , obtained during normative
system operations. The database is meant to be populated before using Au-
toLog; moreover, numeric scores computed from the normative chunks in
the database are used to train the anomaly detector, i.e., *training phase* in
376 Figure 2. It is worth noting that log analysis is a “moving” target: software
upgrades or changes to the configurations may alter meaning and character
of the logs during the lifetime of a given system (Oliner and Stearley, 2007).
While the implementation of a *re-training* component is not within the per-
380 spectives of the paper at this time, the architecture presented in Figure 2
is suited to keep up with changes in the logs. When needed, the support
database of AutoLog can be augmented or updated by operations engineers
with new batches of normative chunks; in turn, the anomaly detector is re-
384 trained with the scores computed from the newly-added chunks. Updating
the database has no specific impact on the overall functioning of our method
if not the time taken to insert the new chunks and training the detector. The
388 description of the *scoring* component is addressed in Section 3.2.2, while the
anomaly detection approach is detailed in Section 4.

¹³<https://logging.apache.org/log4j/>

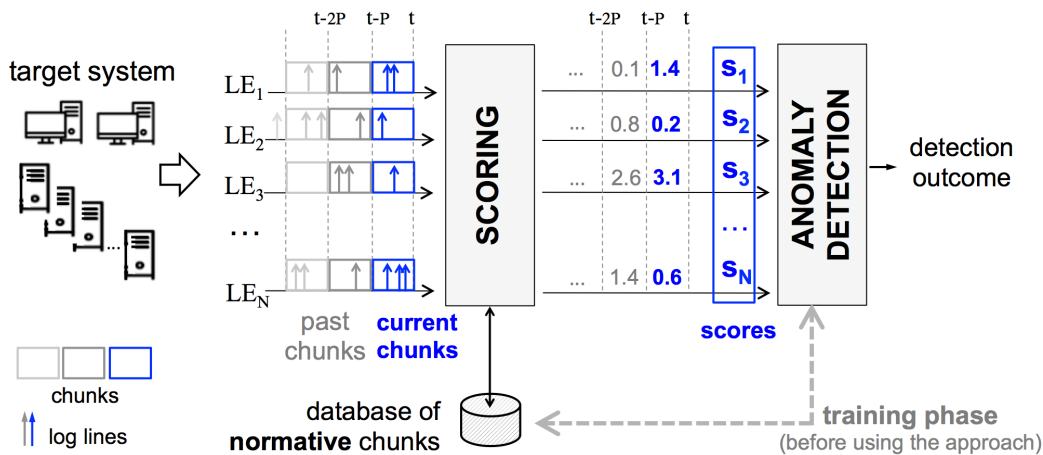


Figure 2: Overview of AutoLog.

3.2.2. Scoring Component

Let LE_i (with $1 \leq i \leq N$) denote a Logging Entity (LE), i.e., a component of the target system (such as a daemon, an application, a node or a set of nodes) that emits lines in the logs. A system may “natively” store the lines emitted by the entities into distinct log files, e.g., one file per application, microservice or container; alternatively, lines from all the entities are centralized at a unique location, and can be dissected later on by source, such as in BG/L.

Regardless how files are physically organized and stored, each LE produces its own timeline of log lines, which are denoted by \uparrow in Figure 2. The **scoring component** acquires the log lines emitted by the LEs in the form of batches of fixed-length windows –named **chunks**– of period P . It is worth noting that the approach we use for handling logs is in-line with the well-consolidated *micro-batch* stream processing pattern adopted by up-to-date massive data processing frameworks, such as Spark streaming¹⁴. At a given time, the scoring component processes the chunks belonging to the same sampling period, i.e., *current chunks* in Figure 2; once processed, current chunks are queued to the *past chunks* and scoring moves on to the next period.

Current chunks are fed to the scoring component represented in Figure 2, which applies *parsing* and *term weighting* to the chunks in order to compute a vector of numeric scores s_i ($1 \leq i \leq N$), i.e., one per LE_i ; it is worth noting that

¹⁴<https://spark.apache.org/streaming/>

a new vector of scores is generated for every period P . **Parsing** is a common data preparation step in log analysis (He et al., 2016a) and allows removing the variable tokens of the log lines while preserving the constant parts. At this stage, we also remove special characters and punctuation, such as #, ?, and %. On the other hand, **term weighting** has been successfully used in the past to operate on text logs (Stearley and Oliner, 2008). Given a chunk after parsing, term weighting is done by (i) tokening the log lines of the chunk into terms¹⁵, (ii) counting the occurrences of the terms within the chunk, (iii) computing a numeric score for the chunk based on the occurrences of the terms. Tokenization is applied to all the lines, including those pertaining to negations in the log files.

Let x_t denote the number of occurrences of the **term** t in the chunk of a given LE_i (with $1 \leq t \leq T$, where T is the total number of terms). The score of the chunk is given by:

$$s = \sqrt{\sum_{t=1}^T (e_t \cdot \log_2(1 + x_t))^2} \quad (0 \leq e_t \leq 1) \quad (1)$$

where e_t is the *entropy* of the term t . The **entropy** is computed from both x_t and by counting the occurrences of t in the set of $(M-1)$ normative chunks of LE_i stored in the support database presented in Section 3.2.1. Figure 3 exemplifies the role of the database at supporting the computation of the entropy. For each LE_i , the database contains $(term, count)$ pairs computed from $(M-1)$ chunks collected under normative operations; noteworthy, $(term, count)$ pairs are arranged by originating chunk, each represented by a dotted box and labeled with 1, 2, ..., $(M-1)$ in Figure 3. As previously mentioned, the database is populated before the use of the approach. Once the scoring component is fed with a new chunk –“current” chunk of LE_i in Figure 3– it generates the $(term, count)$ pairs and retrieves the counts of each term from the database, such as for **enable** in Figure 3. That said, the value of the **entropy** e_t is given by:

$$e_t = 1 + \frac{1}{\log_2(M)} \sum_{j=1}^M p_{t,j} \log_2(p_{t,j}) \quad p_{t,j} = \frac{x_{t,j}}{\sum_{j=1}^M x_{t,j}} \quad (2)$$

¹⁵A term is a sequence of characters separated by one (more) whitespace(s).

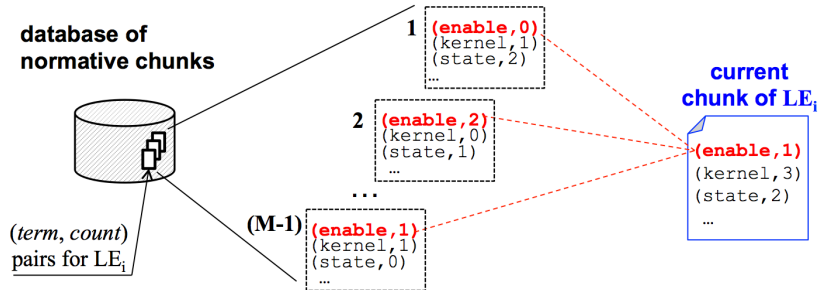


Figure 3: Role of the support database and term-count representation of the chunks used for computing the entropy.

436 where M is the total number of chunks, i.e., $(M-1)$ chunks from the database
plus the chunk targeted by the scoring component (represented in the right-
most part of Figure 3), $x_{t,j}$ is the number of occurrences of the term t in the
chunk j , and $p_{t,j}$ is the fraction of term t 's total occurrences that are in the
440 chunk j (with $1 \leq j \leq M$). We assume that the chunk targeted by the scoring is
the M th chunk: accordingly, $x_{t,M} = x_t$ in Equation 2. Beside the arrangement
into vectors for anomaly detection, each score produced by the scoring compo-
nent is tagged with (i) a unique numeric identifier, which tracks the vector
444 the score belongs, (ii) the originating logging entity, and (iii) time-stamps of
first and last log lines in the chunk used to compute the score. Although not
intended for anomaly detection, administrators can leverage this accessory
data to traverse system logs upon the occurrence of an anomaly.

448 The database of normative chunks serves as a baseline of the system be-
havior. In this respect, the score obtained by applying Equations 1 and 2
quantifies the extent an arbitrary chunk resembles (or not) the baseline of
terms that are emitted by the system under regular operations. The *higher*
452 the score, the *larger* the difference between the terms of the chunk and “typi-
cal” lines emitted by LE_i , which may be seen as an indication of an anomaly.
For example, Table 2 shows some real chunks from BG/L and the correspond-
ing scores computed by our technique after parsing and term weighting. The
456 first chunk is assigned 0.23, i.e., a relatively low value, because it reports
on detected and corrected errors, which are purely informative and can be
noted in the logs of BG/L almost at any time. On the other hand, the latest
two chunks point to real anomalies: they achieve a high score because the
460 log lines consist of terms that are infrequent across the logs. It is worth
noting that the chunk shown by the bottom row of Table 2 results from the

Table 2: Examples of chunks and corresponding scores from BG/L .

log lines in the chunk	score
1 ddr errors(s) detected and corrected on rank 0, symbol 18, bit 6 CE sym 18, at 0x0deb5a60, mask 0x02 total of 1 ddr error(s) detected and corrected	0.23
ddr: Unable to steer rank=0, symbol=5 - rank is already steering symbol 4. Due to multiple symbols being over the correctable error threshold, consider replacing the card1	3.68
machine check interrupt instruction address: 0x001544bc ... omitted ... rts panic! - stopping execution	6.76

interleaving of different anomalous lines, whose “aggregated” effect is a score higher than the individual lines themselves.

464 The weighting approach used here is also known as **logarithmic entropy**
because –according to Equation 1– the dominance of the terms is mitigated
by \log_2 and then scaled by the their entropy. The application of information
entropy as a measure for the uncertainty in a data set and detection purposes
468 is well consolidated (Holzinger et al., 2014). Moreover, there is an extensive
body of literature on term weighting and related applications. Interested
readers are referred to introductory references, such as (Salton and Buckley,
1988; Quan et al., 2011), for additional information on the topic.

472 3.3. Available Datasets

Vectors of scores are computed from the logs reflecting normative operations and anomalies of the systems presented in Section 3.1. In this respect, experiments are based on both simulated and spontaneous anomalies.

476 3.3.1. System Logs

Industrial and microservices systems. System logs from normative (NORM) operations are obtained by exercising the systems with the client applications presented in Section 3.1. In addition, we collect the logs from
480 five independent scenarios, each reproducing a different anomaly. The occurrence of the anomaly is intertwined with the normative operations; after its beginning, each anomaly stays on up to 30 minutes. We adopt a uniform

model of anomalies in the both the industrial and the microservices system;
484 anomalies are spiked in as follows:

- Authentication (**AUTH**): bruteforce attempts to gain unauthorized access to a system. For both the systems in hand, authentication-related anomalies are elicited by attempting to guess the credentials of a legitimate user. This is achieved by means of remote logins via the frontend applications exposed by the systems to external clients (e.g., the human-machine interface of the industrial system).
488
- Log deletion (**DEL**): deletion of the content of a log. It should be noted that this is a typical action performed by an attacker to cover his/her traces.
492
- Hang (**HANG**): the system becomes unresponsive and no service/output is provided within an acceptable timeframe. For example, in the case of the industrial system, the anomaly is injected by stopping the *database service*. As a result, the remaining applications of the system stay up; however, they are not able to successfully fulfill the client requests.
496
- Modification (**MOD**): tampering with the data structures handled by the system through abnormal alterations of their fields. For the industrial system, modification is implemented by changing the expected routes of the vehicles; similarly, in the microservices system we alter the telephone account records.
500
504
- Denial of Service (**DOS**): abnormal usage of the system functions over the nominal capacity with the aim of disrupting service. In both systems, this is achieved by increasing the number of client applications, and thus frequency and volume of the service requests.
508

Although reproduced, anomalies are inspired by the well-consolidated attack phases reported in (Ruiu, 1999); moreover, they cover a mixture of diverse scenarios that range from bruteforce authentication attempts to tampering with OS resources and misuse of system functions.
512

Supercomputing system. It is standard practice to log messages in a supercomputing system, such as BG/L. Log lines in the BG/L log account for hardware and software errors at all levels. Different from the systems

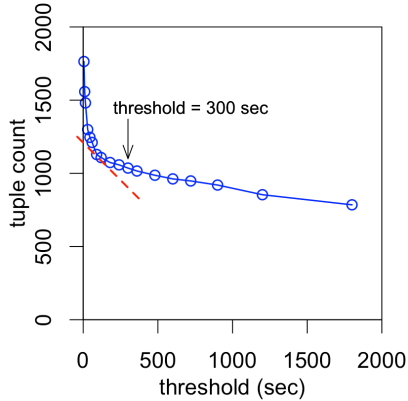


Figure 4: Tuple count of anomalous log lines in BG/L.

516 presented above, the BG/L log contains spontaneous –neither induced or
simulated– anomalies observed over 215 days of operations at LLNL (Oliner
and Stearley, 2007).

Hadoop system. Hadoop logs used in this study are collected dur-
520 ing different executions of data processing applications distributed across a
cluster of machines. Logs pertain to both normative executions of the appli-
cations and different types of service failures in the production environment,
such as *machine down*, *network disconnection* and *disk full*. The dataset is
524 publicly-available as a benchmark for log-based anomaly detection; details
can be found in the proposing paper (Lin et al., 2016).

3.3.2. Sampling and labelling

System logs are sampled into chunks in order to compute the scores for
528 anomaly detection. We set a period $P=10$ s for the industrial, microservices
and Hadoop system, which is a balanced trade-off between the latency of
the detection and the need for ensuring a suitable number of lines per chunk
after each sampling round. As for BG/L, we conduct a sensitivity analysis
532 beforehand. This stems from the long-standing observation that anomalies
tend to cause multiple and redundant log lines, which should be grouped
together. A common technique to address this issue is the **tuple heuristic**
(Hansen and Siewiorek, 1992), which groups the lines whose time distance is
536 lower than a threshold into the same “tuple”. The threshold is determined
by testing different time values and counting the resulting number of tuples
(i.e., the tuple count). Figure 4 shows how the tuple count of the anomalous
log lines of BG/L varies with respect to the threshold. Experimental studies

Table 3: Number of unique n -grams, logging entities (LE), sampling period and chunks by system.

system	corpus statistics			LE	period	number of chunks
	unigrams	bigrams	trigrams			
industrial	898029	2382190	7235966	16	10 <i>s</i>	22640
microservices	2370889	5129706	8689044	13	10 <i>s</i>	12116
BG/L	5632912	12207650	21771832	67	300 <i>s</i>	3473548
Hadoop	152068	328295	541476	28	10 <i>s</i>	114884

540 demonstrate that a good choice for the threshold is the value right after the
“knee” of the curve –dotted line in Figure 4– where the tuple count flat-
tens sharply (Hansen and Siewiorek, 1992). Based on the result, we assume
 $P=300$ *s* for BG/L.

544 Table 3 shows the number of unique unigrams, bigrams and trigrams (thus
providing a glimpse of each log corpus), logging entities, sampling period
and chunks by system. For the industrial, microservices and Hadoop system,
logging entities relate to distinct log files; as for BG/L, entities consist of
548 the set of nodes allotted to the same rack. As a remark, AutoLog produces
one score per chunk: in consequence, –given a dataset– the number of scores
is equal to the number of chunks shown by the rightmost column of Table
3. The number of chunks is determined by the number of logging entities,
552 sampling period and duration of the logs in hand.

Figure 5 shows the scores of two logging entities, i.e., LE_3 for the in-
dustrial system (Figure 5a) and LE_6 for the microservices system (Figure
5b), under normative operations and the occurrence of an anomaly over 30
556 subsequent chunks (i.e., five minutes). It can be noted that the anomalies,
i.e., *MOD* and *DOS* (Δ -marked series), make the score to considerably de-
viate from the normative ones (\circ -marked series). It is worth noting that
normative scores are computed from the logs elicited through representative
560 client applications and benchmarks, which generate mixtures of interleaving
requests. In this respect, the normative scores reflect the variability of the
workload at the time logs were collected. Most notably, the y-axis of Figure
5 is in log scale: while it allows appreciating fluctuations around low values
564 of the scores, their variability is over-emphasized.

Scores are arranged into vectors according to the approach in Section
3.2.2. Furthermore, we accompany each vector by a **label**, which denotes
whether the vector relates to normative or anomalous conditions. For the

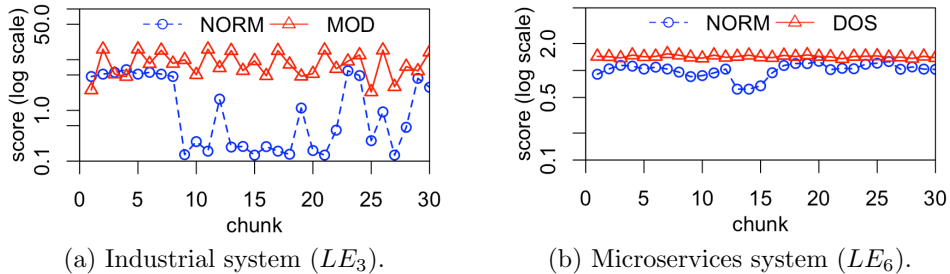


Figure 5: Scores of two logging entities within normative and anomalous conditions.

568 industrial and microservices systems the label is either *NORM* (normal)
or any of *AUTH*, *DEL*, *HANG*, *MOD*, *DOS*. As for BG/L and Hadoop, it
must be noted that the public logs used in our study are labeled. In fact,
logs were tagged with a label established by domain experts in consultation
572 with the system administrators: we rely on this information to label the
vectors produced by our approach, i.e., *NORM* or *ANOM* (anomaly). Whilst
AutoLog *does not* need the labels at training time, labels are intended to be
used for evaluating the effectiveness of both AutoLog and other techniques
576 assessed in the comparative study.

4. Anomaly Detection Method

4.1. Background

AutoLog hinges on the use of a **deep autoencoder**. An autoencoder
580 (AE) is a feedforward neural network where the output layer has the same
dimension of the input layer. In fact, the purpose of an AE is to “reconstruct”
the input at the output layer. It is possible to design different types of
autoencoders (Goodfellow et al., 2016). In particular, deep learning can be
584 applied to autoencoders: multiple hidden layers are used to provide depth.
The resulting network is known as *deep* or *stacked autoencoder* (Vincent
et al., 2010).

An autoencoder learns the input representation in a different feature space
588 (Goodfellow et al., 2016). The learning task forces the autoencoder to catch
the most relevant features of the training data at the **bottleneck** layer,
i.e., the middle hidden layer, so that the input can be reconstructed at the
output layer. Therefore, the AE consists of two parts: *encoder* and *decoder*.

592 The encoder learns an efficient representation of the given input by using
its hidden layer(s); the decoder, instead, produces the reconstruction of the
input by using the encoded information in its hidden layer(s). Therefore, the
decoder mirrors the encoders in the number of hidden layers and neurons.

596 An **encoder** is a deterministic map f_θ that transforms an input vector
 s into its representation y , where $\theta = \{W, b\}$, W is the weight matrix and
 b is the bias vector. On the other hand, a **decoder** is the g_θ function that
maps backwards the representation y to the output z , i.e., the reconstruction
600 of the input vector s . In Equation 3 the variables W and b represent the
weight and bias values for the encoding phase. Similarly, in Equation 4 the
variables W' and b' are the weight and the bias for the decoding phase.

$$f_\theta(Ws + b) = y \quad (3)$$

$$g_\theta(W'y + b') = z \quad (4)$$

The **reconstruction error (RE)** measures the difference between the
604 reconstructed, i.e., z , and the original version of the input, i.e., s . It is
important to feed data to an autoencoder and tune it until it is well trained
to reconstruct the input with minimum error. Therefore, the measure of
faithful reproduction of input data is defined by the RE, which is computed
608 as follows:

$$RE = \frac{1}{N} \sum_{i=1}^N (z_i - s_i)^2 \quad (5)$$

where z_i and s_i (with $1 \leq i \leq N$) denote the components of the output and input
vector, and N is the number of components. Overall, the key characteristics
of autoencoders are:

- 612 • *Data dependency.* Autoencoders will only be able to reconstruct data
similar to that on which they have been trained. This principle is the
basis of our anomaly detection approach (more on this later).
- *Output lossy.* This means that the reconstructed output will be de-
616 graded compared to the original input.
- *Self-contained learning.* Autoencoders are trained automatically from
data examples. It is not necessary to do the extra work of preparing
extra labels or data.

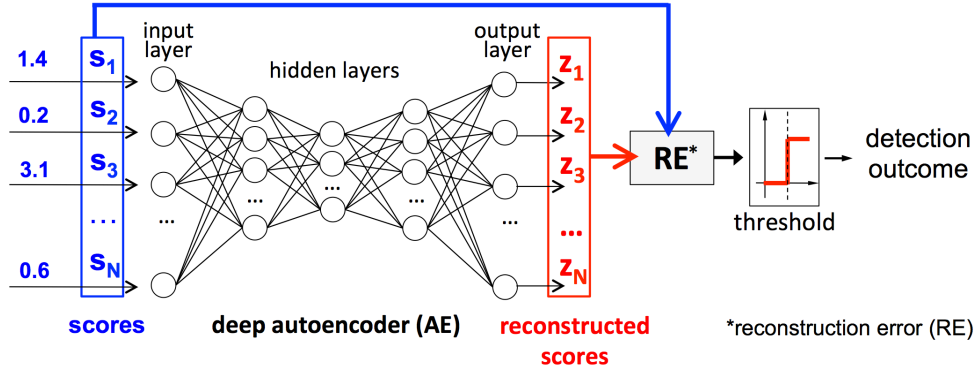


Figure 6: Anomaly detection in AutoLog.

620 An autoencoder can be used for dimensionality reduction, in a way similar
to Principal Component Analysis (PCA) (Almotiri et al., 2017). It is worth
noting that PCA uses linear algebra to make the transformation of the coordi-
624 nates; in contrast, an autoencoder performs non-linear transformations by
means of non-linear activation functions and multiple layers. Hence the use
of an autoencoder is particularly attractive when the data points are complex
and non-linear in nature (Song et al., 2013). This is the case of the datasets
in hand, as it will be shown in Section 5.1.

628 4.2. Use of the Deep Autoencoding in AutoLog

Figure 6 shows a representation of the AE in the context of AutoLog
and it reproduces input and output of the *anomaly detection* block in Fig-
ure 2 (i.e., *scores* and *detection outcome*, respectively). More importantly,
632 Figure 6 details the internal organization of the anomaly detector. Scores
–generated by the scoring component– are fed at the input layer of the AE:
scores pass through a number of hidden layers until the output layer returns
the reconstructed scores. Noteworthy, input and output layers have the same
636 dimension N . In the context of our study, the value of N is selected equal
to the number of logging entities, i.e., components emitting lines in logs as
described above. The number of logging entities and their selection for each
system is presented in Section 3.3.2.

640 The rationale underlying the use of the AE in AutoLog is that RE can
be used as an *anomaly indicator* as follows. If the AE is trained using only
vectors of normative scores –hence the notion of **semi-supervised learn-**
ing– it will provide (i) *low* RE (good reconstructed representation) for future

644 normative input vectors, and (ii) *high* RE (bad reconstructed representation)
for future anomalous input vectors. In this respect, the AE serves as a nor-
mative baseline: when the AE attempts to reconstruct a data point that
is “outside” the norm –and thus a potential anomaly– it will experience an
648 increase of the RE because it was never trained to reproduce anomalies.

As for many anomaly detection techniques, we apply a binary *threshold*
function to RE in order to pinpoint anomalies: the input vectors that pro-
duce RE values under the threshold are deemed *normative* and those with
652 REs above the threshold, *anomalous*. The threshold is called here **anomaly**
threshold: as the specific configuration of the AE for the data in hand (e.g.,
in terms of layers, neurons and activation functions), the setup of a suitable
threshold is an outcome of the training phase and it is addressed in Section
656 4.4. RE computation and comparison with the threshold represent the steps
of AutoLog that produce the final detection outcome, as shown in Figure 6.

4.3. Dataset partitioning

As discussed in Section 3.3, we conduct our experiments with datasets of
660 scores arranged into labeled vectors; labels indicate whether a given vector
is collected from normative operations or under an anomaly. As for the
labels, we carry out a symbolic-numeric conversion by substituting 0 to the
labels corresponding to a normative vector and 1 to the labels corresponding
664 to anomalous vectors. Each dataset is split into three **disjoint subsets**
through random sampling. It is worth noting that the sampling procedure
is *without replacement*: once selected from a set, the vector is not placed
back to the set it comes from. As said, a vector consists of chunk-wise
668 scores computed from the log lines emitted by the logging entities during
a sampling round. Figure 7 shows how a given dataset is partitioned. Let
 NV be the cardinality of the normative vectors and AV the cardinality of
the anomaly vectors. According to Figure 7, the first cut consists in (i)
672 separating normative from anomalous vectors, and (ii) splitting normative
vectors into two disjoint subsets of cardinalities $(0.8 \cdot NV)$ and $(0.2 \cdot NV)$,
respectively. At the bottom of Figure 7, we find:

- 676 • *Training set*. It contains only normative vectors: the cardinality of the
training set is $0.9 \cdot (0.8 \cdot NV)$ of randomly selected normative vectors
from the originating dataset. The set is meant for training the AE.
Labels are removed.

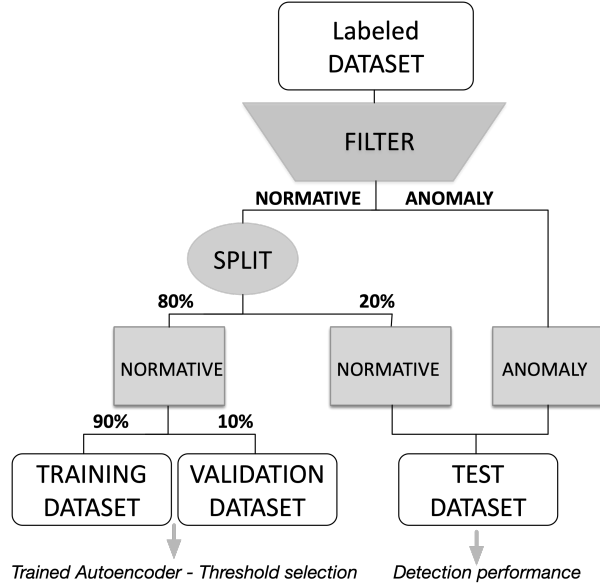


Figure 7: Dataset partitioning

- 680

Validation set. Again, it contains only normative vectors and is used for determining the detection threshold. In particular, the cardinality of the validation set is $0.1 \cdot (0.8 \cdot NV)$. As with the training set, labels are removed.
- 684

Test set. It contains both normative and anomalous vectors. The cardinality of the test set is $(0.2 \cdot NV) + AV$, i.e., 20% normative vectors –obtained at the first cut– plus all the anomalous vectors. According to the sampling procedure and the partitioning in Figure 7, normative vectors of the test set are **held out** from training. Vectors in the test set are accompanied by the corresponding labels in order to evaluate the correctness of the predictions.

692
 After partitioning, we obtain three non-intersecting subsets. Table 4 provides the size of each dataset for the systems in hand, and the corresponding breakdown into training, validation and test set. Given a system, the number of vectors is equal to the number of chunks divided by the number of logging entities of that system. For example, the total chunks of the industrial system (i.e., 22640) returns $\frac{22640}{16} = 1415$ total vectors, where 16 is the number

Table 4: Training, validation and test set size.

system	total chunks/vectors	breakdown		
		training	validation	test
industrial	22640/1415	6928/433	768/48	14944/934
microservices	12116/932	6552/504	728/56	4836/372
BG/L	3473548/51844	2450525/36575	272288/4064	750735/11205
Hadoop	114884/4103	70756/2527	7868/281	36260/1295

Table 5: Layering structure of different AE configurations (all the layers are dense).

Configuration 1		Configuration 2		Configuration 3	
layer	activation	layer	activation	layer	activation
Input	-	Input	-	Input	-
Hidden 1	ReLU	Hidden 1	tanh	Hidden 1	tanh
Hidden 2	tanh	Hidden 2	sigmoid	Hidden 2	tanh
Hidden 3	ReLU	Hidden 3	sigmoid	Hidden 3	tanh
Output	ReLU	Output	tanh	Hidden 4	ReLU
				Hidden 5	sigmoid
				Output	ReLU

696 of logging entities shown in Table 3. A similar computation applies to all
chunks/vectors pairs in Table 4.

4.4. AE Design, Training and Threshold Selection

700 The design of a deep neural network, such as the AE, is based on estab-
lishing many hyperparameters that are subject to fine-tuning. As for any
machine learning study, the choice of the hyperparameters is guided by ex-
perimental tests carried out by analyzing the outcome of the model –RE in
704 our study– with respect to the **validation set**.

There are two *desirable* properties of the RE, which make it possible to
usefully deploy an AE for anomaly detection in our context (i) $RE \approx 0$ and (ii)
small dispersion. An AE design that meets these properties is summarized
708 by the Configuration 1 in Table 5, which is the selected configuration for
the datasets in hand: its RE on the validation set of the industrial system
is shown in Figure 8a, where it can be noted that it is ≈ 0 for all –if not
one– data points. In order to provide concrete examples of less effective AE

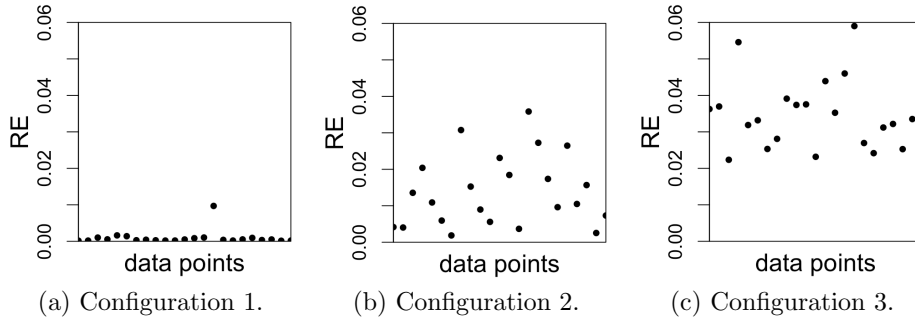


Figure 8: Analysis of RE on the validation set by configuration; Configuration 1 is the final selection for AutoLog.

712 designs for the industrial system, Figure 8b and 8c show the REs achieved
 on the validation set by two different configurations that fail to obtain $RE \approx 0$
 (Configuration 2 and 3 in Table 5).

As summarized by Configuration 1 in Table 5, the chosen AE is made
 716 up of five layers. These layers include N -128-64-128- N neurons, where N
 is the number of logging entities. Dropout layers are placed after the first
 hidden layer and after the second hidden layer, in order to prevent overfitting.
 The Rectified Linear Unit (ReLU) has been selected for the encode layer,
 720 the decode layer and the output layer, while for the bottleneck layer, i.e.,
 Hidden 2, has been used the Hyperbolic Tangent (Tanh) activation function.
 Moreover, to achieve the sparsity, activity regularizer terms L1 are applied on
 each layer. We train the AE on training data points for 100 epochs using the
 724 RMSProp optimizer with learning rate value $lr=0.001$. We also shuffle the
 training data before each epoch. It is worth pointing out that the training
 phase takes around 1 minute in the worst case (BG/L dataset). This time has
 been obtained on a laptop computer without GPU acceleration, and so it can
 728 be greatly reduced using more powerful or ad hoc hardware. In particular,
 the experiments are conducted on a MacBook Pro with an Intel Core i5 2.6
 GHz processor and 8 GB of RAM.

We select the threshold value by using normative points, and hence the
 732 training and validation procedure is semi-supervised. As no label is required,
 this is by far the simplest approach for threshold selection. In particular,
 we adopt a *percentile* method and choose as **anomaly threshold** the 90th
 percentile of the RE values obtained by the AE on the validation set. Section
 736 5.3 validates the proposed threshold selection method by analyzing receiver

operating characteristic curves. Selecting the 90th percentile appears suitable for all the datasets we have considered so far. In practice, it aims to mitigate the impact of sporadic RE outliers caused by accidental anomalies in the normative logs. We are aware that selecting the 90th percentile is not a *one-fits-all* approach. For example, in those systems where anomalies are sparse and with a high false positive rate, the detection approach might further benefit from the adoption of a more sophisticated threshold selection, such as the one recently proposed in (Carrington et al., 2021). Nevertheless, it should be noted that our choice is inline with many other studies in the area that rely on easy-to-explain thresholds, such as (Aygun and Yavuz, 2017), which adopts the mean value. Although worthy to be investigated, digging into fine-grained threshold-related aspects is beyond the perspective of the study at this stage, where we focus on the overall methodology.

We implemented our architecture with Keras¹⁶ (Version 2.4.3) and TensorFlow¹⁷ (Version 2.4.1). Keras is a Python library that runs on top of TensorFlow; it provides highly modularized APIs for building and training deep learning models. The core code of AutoLog encompassing the implementation of the autoencoder has been made publicly-available through GitHub¹⁸.

5. Experimental Results

AutoLog is applied to the reference systems –industrial, microservices, BG/L and Hadoop– in order to quantify its effectiveness to discriminate normative operations from the occurrence of anomalies. In the following, we present some reflections on the challenges of PCA and clustering and discuss the results of AutoLog. A comparison study of AutoLog with a wide set of techniques is provided in Section 6.

5.1. Reflections on the Use of Clustering

We conduct an **exploratory data analysis** to check whether the scores obtained in face of normative conditions and anomalies tend to cluster into different groups. To this aim, the vectors of the datasets can be conveniently regarded as “points” of a Euclidean space. Nevertheless, the *dimensionality* of the datasets, i.e., 16, 13, 67 and 28 –not including the label– for each

¹⁶<https://keras.io/>

¹⁷<https://www.tensorflow.org/>

¹⁸<https://github.com/ScalingLab/AutoLog>

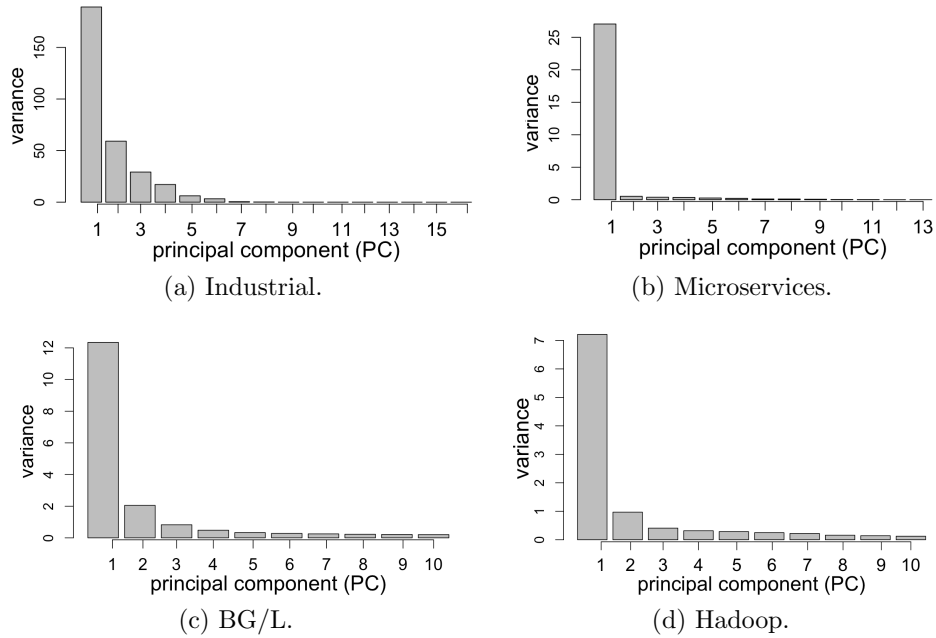


Figure 9: Screeplots of the variance of the principal components.

768 system, respectively (with each dimension corresponding to the number of
 logging entities, as in Table 3) prevents from obtaining human-readable visu-
 alizations for exploratory purposes. In consequence, a Principal Component
 Analysis (PCA) is done beforehand.

772 PCA is a **dimensionality reduction** technique whose objective is to find
 the *directions* along which a set of high-dimensional points line up “best”.
 Points are transformed in a new coordinate space where each axis catches
 progressively less variance of the original data. This concept is summarized
 776 by means of a *screeplot*, which is one of the byproducts of a PCA. Figure 9
 shows the screeplots obtained for the datasets in hand. The x-axis refers to
 the principal components (PC) –each denoted by an integer– and the y-axis
 is the variance explained by the PC; it should be noted that PCs are sorted
 780 by descending variance. For the industrial system (Figure 9a) the variance
 of the top-3 PCs accounts for 189.2, 59.1 and 29.2, respectively, out of the
 total variance of 304.9 (i.e., the sum across all PCs): as a consequence, the
 top-3 PCs catch the 91.0% of the total variance. Similarly, the total variance
 784 of the PCs for the microservices dataset is 29.1; the top-3 PCs account for

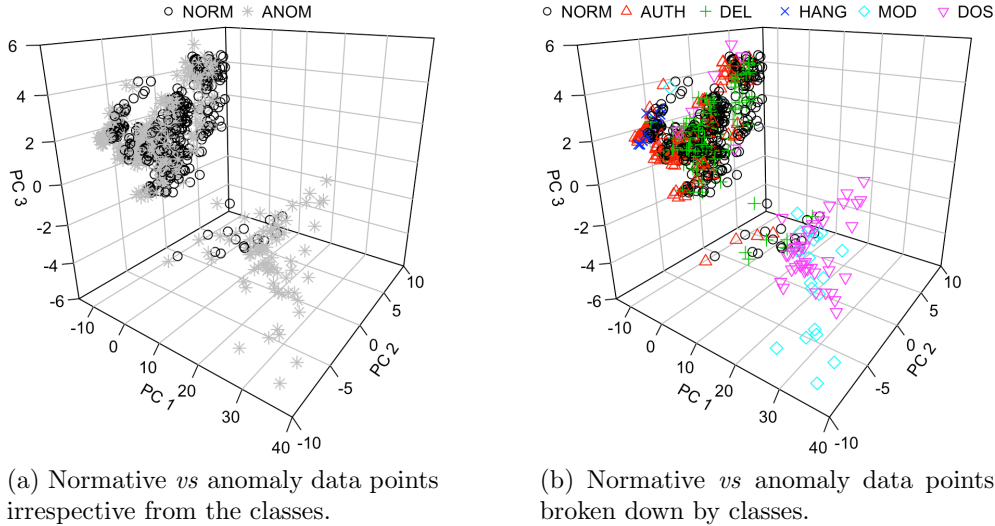
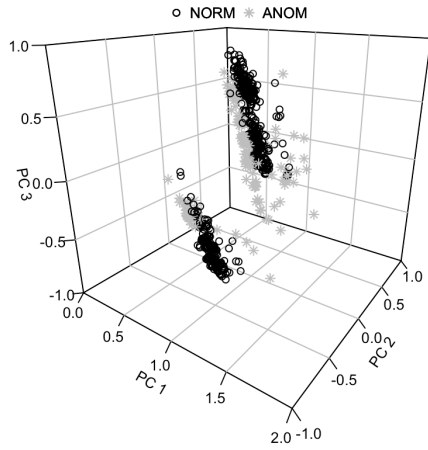


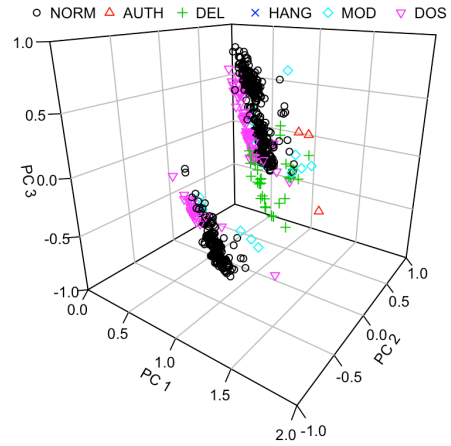
Figure 10: Scatterplot of the dataset with respect to the top-3 principal components (industrial system).

27.0, 0.52 and 0.39 variance, which sum up to 95.9% of the total. As for BG/L and Hadoop in Figure 9c and 9d, we show the top-10 components due to the large number of original dimensions: the top-3 PCs encompass 68% and 80% of the total variance for BG/L and Hadoop, respectively. Based on these results, we *narrow* the visualization to a more suitable 3D Euclidean space, where each point consists of the coordinates along the top-3 PCs.

Figure 10a shows the 3D scatterplot of the industrial system dataset, where the axes are PC 1, PC 2 and PC 3, respectively. Moreover, \circ -marked points denote normative points and $*$ -marked points depict anomalies, i.e., any of the types presented in Section 3.3. It can be noted that the points tend to group around two areas. Whilst at the bottom part of the grid we find almost only anomalous points, the top left group is strongly cluttered: in fact, normative and anomalous points are intertwined. Figure 10b provides a similar perspective where anomalies are broken down by class: again, \circ -marked points denote normative points while the remaining mark types are meant to represent anomalies. It can be noted that misuse-related anomalies, i.e., *MOD* (\diamond mark) and *DOS* (∇ mark), take a clear stand from normative points; however, the real challenge for detection is the *intertwinement* of 4 different classes in the top left group, which undermines the practical us-

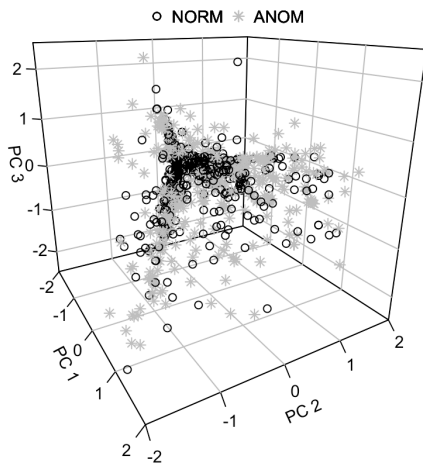


(a) Normative *vs* anomaly data points irrespective from the classes.

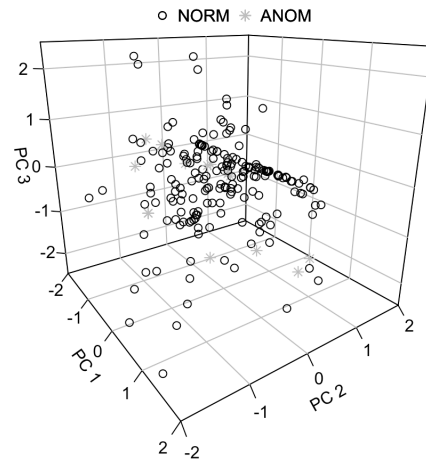


(b) Normative *vs* anomaly data points broken down by classes.

Figure 11: Scatterplot of the dataset with respect to the top-3 principal components (microservices system).



(a) BG/L.



(b) Hadoop.

Figure 12: Scatterplot of BG/L and Hadoop datasets with respect to the top-3 principal components.

804 ability of PCA and clustering with respect to the data in hand. This finding
 holds for the microservices dataset as well. Figure 11 shows the 3D scatterplot
 of the top-3 PCs for the microservices dataset, i.e., *binary* view (Figure 11a)
 and *multiclass* view (Figure 11b). Again, normative and anomalous points
 808 are strongly intertwined.

The scarce *linear separability* of the data points is a key challenge to the
 use of PCA and clustering as reference technique for unsupervised learning.
 While the industrial and microservices systems were available at a private
 812 premise, BG/L and Hadoop are common benchmarks being used by other
 papers in the area of log analysis. Figure 12 indicates that the data points
 of BG/L and Hadoop suffer from the same intertwining of normative and
 anomaly classes. For example, it is worth noting that the paper (Meng et al.,
 816 2019) attempts to apply PCA to BG/L with no success for anomaly detection.
 Our deep autoencoding approach capitalizes on non-linear transformations
 of the data, which are much more suited to the datasets in hand and yield
 to better results.

820 5.2. Evaluation of AutoLog

We run the **test set** of the datasets in hand with AutoLog trained as
 described in Section 4.4. Since the datasets are labeled, each RE produced
 by AutoLog is accompanied by the label, which we use for evaluation; as
 824 said, the AE saw *no anomalies* during training. Figure 13, 14, 15, and 16
 show the REs for the industrial system, the microservices system, BG/L, and
 Hadoop, respectively. In particular, each data point is marked by either \circ or
 ∇ for better visualization of normative and anomalous points; moreover, we
 828 superimpose vertical continue lines to delimit the data points by class, which
 are named according to the descriptions in Section 3.3.2. A semi-logarithmic
 scale (x-axis in linear scale and y-axis in log scale) is used to better visualize
 the RE values. The y-axis is the RE; the x-axis is the id of the points in the
 832 test set. The horizontal dashed line in Figure 13, 14, 15 and 16 indicates the
 anomaly threshold, which is 0.00079, 0.00460, 0.00080, and 0.03937 for each
 system, respectively. Section 5.3 investigates the validity of our threshold
 selection approach.

836 We compute the metrics of *recall* (R), *precision* (P), and *F1 score* to
 evaluate AutoLog. Metrics are computed as:

$$R = \frac{TP}{TP + FN} \quad P = \frac{TP}{TP + FP} \quad F1 \text{ score} = 2 \cdot \frac{P \cdot R}{P + R} \quad (6)$$

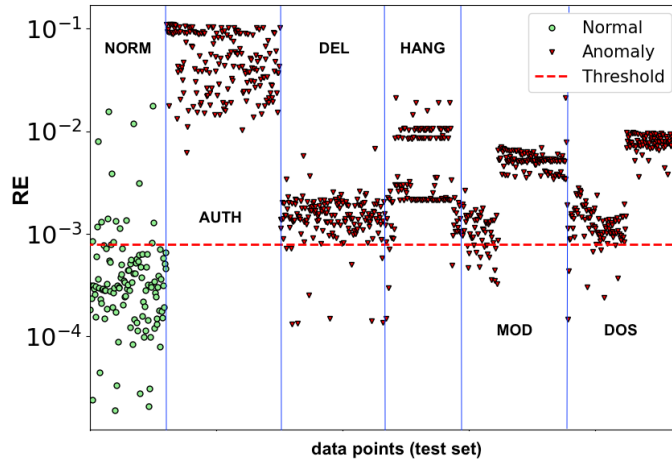


Figure 13: AutoLog: RE of the test set for the industrial system.

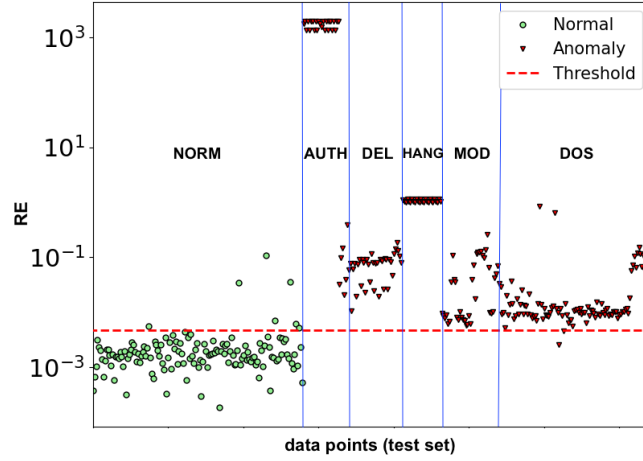


Figure 14: AutoLog: RE of the test set for the microservices system.

840 where True Positive (TP) and True Negative (TN) represent the points that are correctly predicted, while False Positives (FP) and False Negatives (FN) indicate misclassifications. For example, TP is the set of anomalies whose

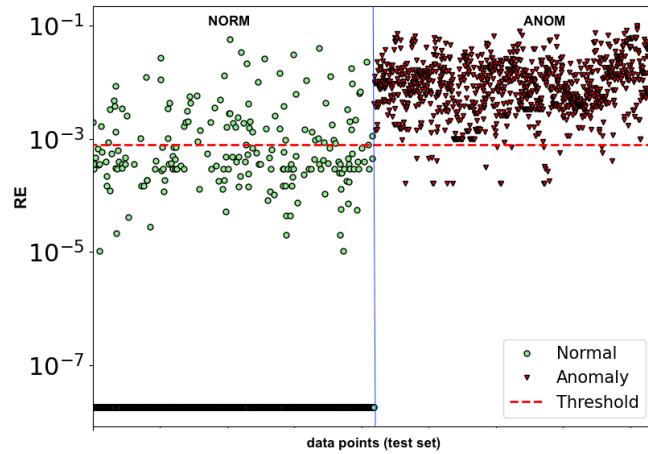


Figure 15: AutoLog: RE of the test set for BG/L.

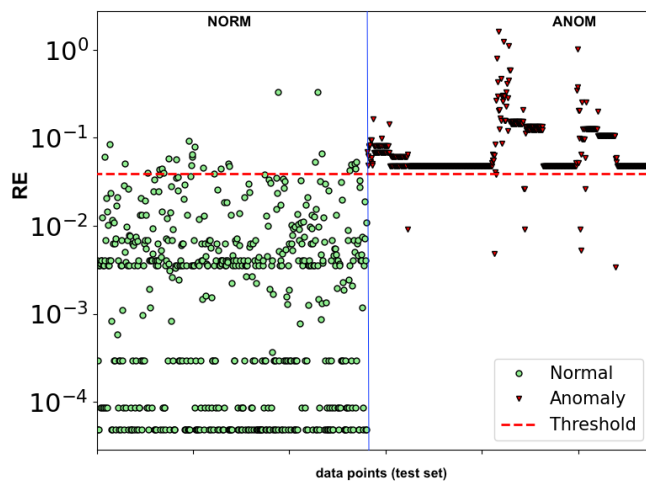


Figure 16: AutoLog: RE of the test set for Hadoop.

844 RE is higher than the threshold; similarly, TN is the set of normative points whose RE is lower than the threshold. Table 6 provides the evaluation metrics of AutoLog for all the systems.

Table 6: AutoLog: evaluation metrics.

system	recall	precision	F1 score
industrial	0.96	0.98	0.97
microservices	0.99	0.97	0.98
BG/L	0.98	0.93	0.95
Hadoop	0.98	0.96	0.97

Results are notable, especially if we consider that AutoLog is based on a semi-supervised approach, which means that it does not need for anomalies at training time. According to Figure 13, 14, 15 and 16, we observe that most of the anomalies have high RE and are well above the thresholds for all the systems. This outcome allows making some relevant considerations. For example, unlike the results obtained by means of PCA, *AUTH*, *DEL* and *HANG* anomalies of the industrial system are now easily identifiable: their RE is much higher than the typical RE of normative data points. These anomalies were strongly intertwined with normative points in the scatterplots of the dataset shown in Figure 10. A similar consideration can be done for the microservices system. According to Figure 11, data points grouped in two clusters –both composed by normative and anomalous points– in the Euclidean space: differently, the points appear well-separable with AutoLog, as in Figure 14. Furthermore, Figure 15 and 16 show that anomalies are above the threshold also for BG/L and Hadoop.

The discussion in Section 4.4 (where we presented the selection of the hyperparameters of the autoencoder and how they affect the RE) is supplemented here by means of the assessment of anomaly detection performance. As for Configuration 2 and 3 reported in Table 5 –strongly different from AutoLog in terms of both activation functions (Configuration 2), and number of layers and activation functions (Configuration 3)– they achieve R=0.89, P=0.97, F1 score=0.93 and R=0.20, P=0.85 and F1 score=0.32, respectively. These detection figures are unsatisfactory when compared to AutoLog reported in Table 6: both the configurations fail to render RE \approx 0 on the validation set –one of the key design principles of AutoLog– as shown in Section 4.4.

We perform an additional analysis by assessing the sensitivity of the F1 score while *narrowing*, *widening* and *deepening* the AutoLog autoencoder. Results are summarized in Figure 17 for the industrial system dataset. As

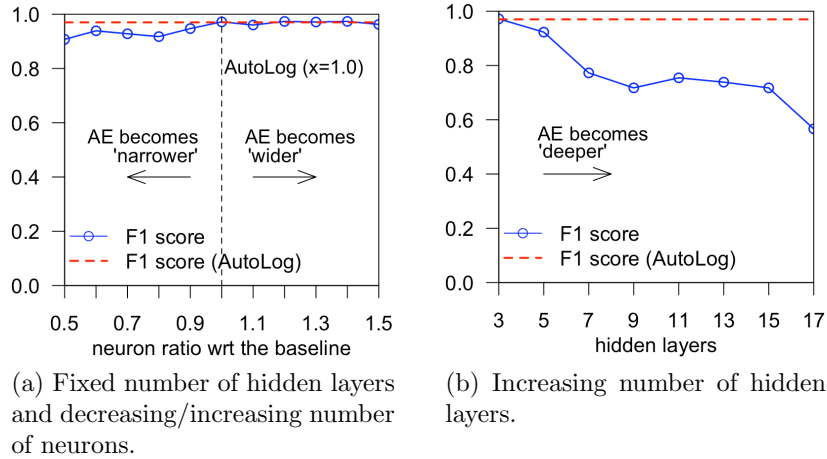


Figure 17: Sensitivity of the F1 score with respect to (wrt) the baseline AutoLog autoencoder (AE).

for Figure 17a, we test different autoencoders with 3 hidden layers. The x-axis is the **neuron ratio** with respect to (wrt) the selected number of neurons for AutoLog at the hidden layers, i.e., (128,64,128): given a value of x in Figure 17a, the number of neurons of the autoencoder under test is given by (128,64,128) $\cdot x$. For example, $x=0.5$ returns (64,32,64) neurons, i.e., a “narrower” autoencoder compared to AutoLog; similarly, at $x=1.5$ we obtain (192,96,192), i.e., a “wider” autoencoder. It is worth noting that $x=1.0$ corresponds to AutoLog itself, whose F1 score for the industrial system is represented by a dashed horizontal line. Figure 17a indicates that the F1 score achieved by the autoencoders flattens at $x=1.0$, i.e., the number of neurons we selected for AutoLog; *widening* further the autoencoder does not improve over AutoLog. On the other hand, Figure 17b shows how the F1 score varies when we increase the number of hidden layers. As indicated by the x-axis, analysis is done by step 2 because we inject two hidden layers –one at the encoder and one at the decoder, according to Section 4– for each test; moreover, the number of neurons for the new layers is set between 64 and 128. Noteworthy, the leftmost data point, i.e., 3 hidden layers, corresponds to AutoLog. Figure 17b indicates that, with respect to the problem addressed and data in hand, deepening the autoencoder does not necessarily improve anomaly detection.

5.3. Validation of the Anomaly Threshold

Choosing the optimal anomaly threshold is a critical task, especially if
896 performed automatically (Liu et al., 2015). As explained in Section 4, our
threshold selection is based only on the REs of normative data points in the
validation set; notwithstanding, the threshold is able to discriminate fairly
well normative from anomalous points. In order to validate our threshold se-
900 lection method, we compare it with the “best-fit” threshold selection that can
be done on the test sets when considering both the RE and the *ground truth*,
i.e., the knowledge of the class of the points to be classified. The availability
of the ground truth is the typical assumption of supervised techniques.

904 We analyze the Receiver Operating Characteristic (**ROC**) curve for our
model. It can be considered as a diagnostic plot, which evaluates the effec-
tiveness of the classification at various threshold settings. The ROC curve
shows how much the model is capable of distinguishing between classes at
908 different operating points. The area under the ROC curve (**AUC**) provides
a numeric score to summarize the performance of a model with a value be-
tween 0.5 (no-skill) and 1.0 (perfect skill). The higher the AUC value, the
better the model at predicting the classes. The ROC curve is plotted with
912 **True Positive Rate** (y-axis) against the **False Positive Rate** (x-axis). A
diagonal line on the plot from the bottom-left to top-right indicates the curve
for a *no-skill* model, and a point in the top left area of the plot indicates a
model with *perfect skill*. The performance of each classifier is represented by
916 a point of the ROC curve and, as the threshold changes, the location of the
point changes. The point that allows the best performance is in the top left
area of the plot. From the analysis of the ROC curves corresponding to our
model for the datasets in hand (Figure 18) it is possible to note a number of
920 points –and therefore potential thresholds– in the top-left area of the plots.

Our goal is to locate the threshold with the optimal balance between
false positive and true positive rates. There are several techniques that allow
to achieve this goal (Du et al., 2017a). In particular, for our analysis we
924 have selected the **Geometric Mean or G-Mean** score. It is a metric that
tries to find a balance between the **Sensitivity** and the **Specificity**, and
computed as follows:

$$G\text{-Mean} = \text{sqrt}(\text{Sensitivity} \cdot \text{Specificity}) \quad (7)$$

where:

$$\text{Sensitivity} = \text{TruePositiveRate} \quad (8)$$

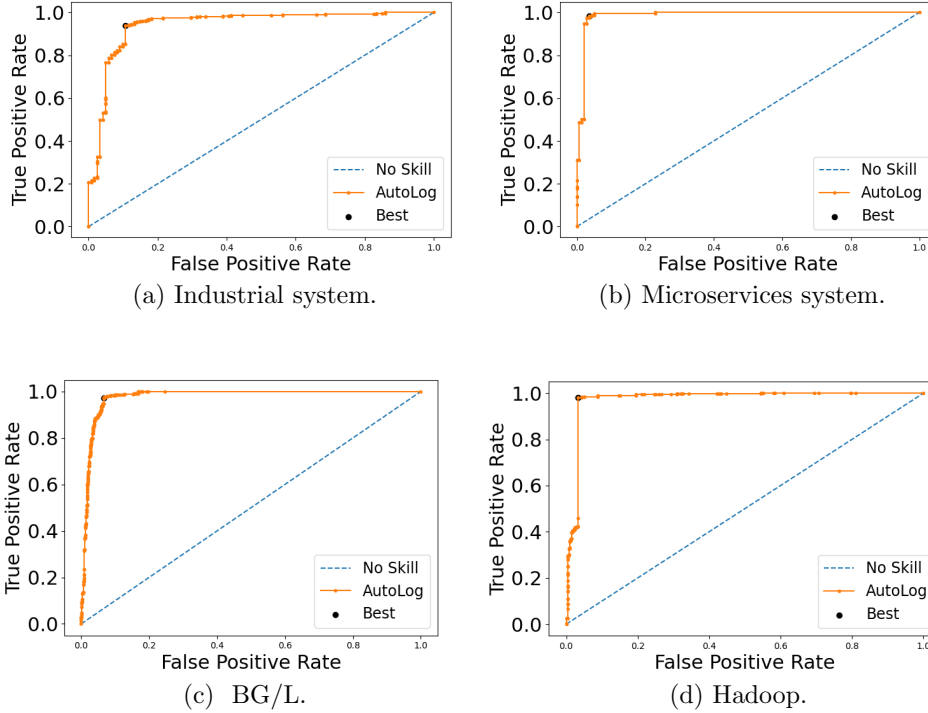


Figure 18: AutoLog: ROC curves

928

$$Specificity = (1 - FalsePositiveRate) \quad (9)$$

The approach we have pursued is to assess all the thresholds obtained from the ROC analysis and select the one with the highest G-Mean score. The ROC curves in Figure 18a, 18b, 18c and 18d highlight the optimal points (●-marked points in the top left area of the plots). For a direct evaluation we report in Table 7 the thresholds produced by the proposed percentile selection method and the best-fit thresholds obtained with the ROC curve method for all the datasets. It is worth noting that the threshold values obtained with the two methods are very close. The two approaches lead to very similar results, and so that our choice –semi-supervised and based on a much simpler procedure– is perfectly reasonable.

932

936

Table 7: Anomaly threshold values.

system	percentile method (AutoLog)	ROC curve method (supervised scenario)
industrial	0.00079	0.00086
microservices	0.00460	0.00564
BG/L	0.00080	0.00100
Hadoop	0.03937	0.04821

6. Comparative Study

940 In this section we compare AutoLog with the following techniques: **iso-**
lation forest, one-class SVM, decision tree, vanilla autoencoder and
variational autoencoder. These techniques are widely used in different re-
research fields for detecting anomalies from data and are applied to the datasets
944 in hand. We implement the isolation forest and one-class classifier by means
of `scikit-learn`¹⁹. The decision tree is implemented with WEKA²⁰; both
the vanilla and variational autoencoder are based on Keras. Given a dataset,
for all the techniques, if not the decision tree, we use the same training and
948 test conditions of AutoLog: (i) 80% normative vectors –random selection
without replacement from the originating dataset– for training and, (ii) re-
remaining 20% normative vectors *plus* all the anomaly vectors for testing and
computing the evaluation metrics. As for the decision tree, we use both
952 normative and anomalous vectors for training and testing, as described in
Section 6.3. In the following sections we briefly introduce each technique
and then present recall, precision and F1 score achieved for each system.

6.1. Isolation Forest

956 Isolation forest (or *iForest*) is an anomaly detection technique based on
the assumption that anomalous data points are rare and far from the centers
of normal clusters (Liu et al., 2008). Many anomaly detection techniques rely
on the construction of *normal profiles* by defining how “normal” points look
960 like: in consequence, anomalies are those points that do not conform to the
defined “normal” profile. Different from this general approach, an isolation
forest aims to directly target anomalies.

¹⁹<https://scikit-learn.org/stable/>

²⁰<https://www.cs.waikato.ac.nz/ml/weka/>

Table 8: Isolation forest: evaluation metrics.

system	recall	precision	F1 score
industrial	0.53	0.87	0.66
microservices	0.85	0.87	0.86
BG/L	0.64	0.78	0.70
Hadoop	0.53	0.71	0.61

The isolation forest technique generates **partitions** on a given dataset by randomly selecting a feature from the given set and a split value between the minimum and maximum values of the selected feature. Anomalous points are likely to require fewer partitions to be isolated compared to the so defined “normal” data points in the dataset. Partitioning is represented by a **tree structure**: anomalies will be the points with a shorter path length within the tree, i.e., number of partitions required to isolate the points. As with other anomaly detection methods, also in the case of isolation forest is defined an *anomaly score*:

$$s(x, n) = 2^{-\frac{E(h(x))}{c(n)}} \quad (10)$$

let x the point and n the number of external nodes (with no child), $h(x)$ is the path length of the point x and $c(n)$ is the average path length of unsuccessful search in the tree. Each point is assigned an anomaly score value.

For our analysis we select the number of *base estimators* –number of trees in the forest– equal to 100. The value is selected after a series of tuning experiments. Table 8 shows the evaluation metrics for the datasets in hand. The best results of the isolation forest are obtained with the microservices system, although none of the metrics is higher than 0.90 in any system.

6.2. One-class SVM

Support Vector Machine (SVM) is a supervised technique frequently used in classification problems (Pisner and Schnyer, 2020). In general, it leverages **hyperplanes** in multi-dimensional space in order to separate classes of points. The distance between the nearest points is known as the *margin*. The core idea of SVM is to find a Maximum Marginal Hyperplane (MMH) that best separates the points into classes. Given labeled training data, the technique produces an optimal hyperplane that allows to classify new points.

Table 9: One-class SVM: evaluation metrics.

system	recall	precision	F1 score
industrial	0.59	0.88	0.71
microservices	0.84	0.87	0.85
BG/L	0.92	0.93	0.92
Hadoop	0.68	0.78	0.73

Although SVM is typically used to solve binary classification problems, it can also be leveraged for **one-class problems**, where the training data belong to one class. In this context, the model is trained to learn what is “normal” so that it can identify whether new data points belong to the “normal” class or not. The technique captures the density of the majority class and classifies data points on the extremes of the density function as outliers. This alternative SVM is referred to as **one-class SVM** (Schölkopf et al., 2001).

In order to apply the one-class SVM technique to the datasets in hand, we leverage `scikit-learn`. The main difference with a standard SVM is that a one-class SVM is fitted in an unsupervised manner and does not provide the typical hyperparameters for tuning the margin. In particular, it provides a hyperparameter nu that controls the sensitivity of the technique and should be tuned to the approximate ratio of outliers in the data. For our experiments we set $nu=0.05$, which we found to be a suitable value for the reference datasets after tuning. Once fitted, the model is used to identify anomalies by classifying the points of the test set as inliers and outliers. Table 9 shows the results obtained for the datasets, where it can be noted that one-class SVM achieves the best results with BG/L (i.e., recall and precision equal to 0.92 and 0.93, respectively).

6.3. Decision Tree

Decision trees are typically used for the capability to infer explicable rules from data. Given an input data point to be classified, decision is made by traversing a tree-like structure starting from the root node. Each node is characterized by a *predicate* meant to be tested on the input data point: based on the outcomes of the tests, decision moves down through the tree until a *leaf* is reached. In a classification problem leaves are associated to the classes the input points are expected to belong. It should be noted that

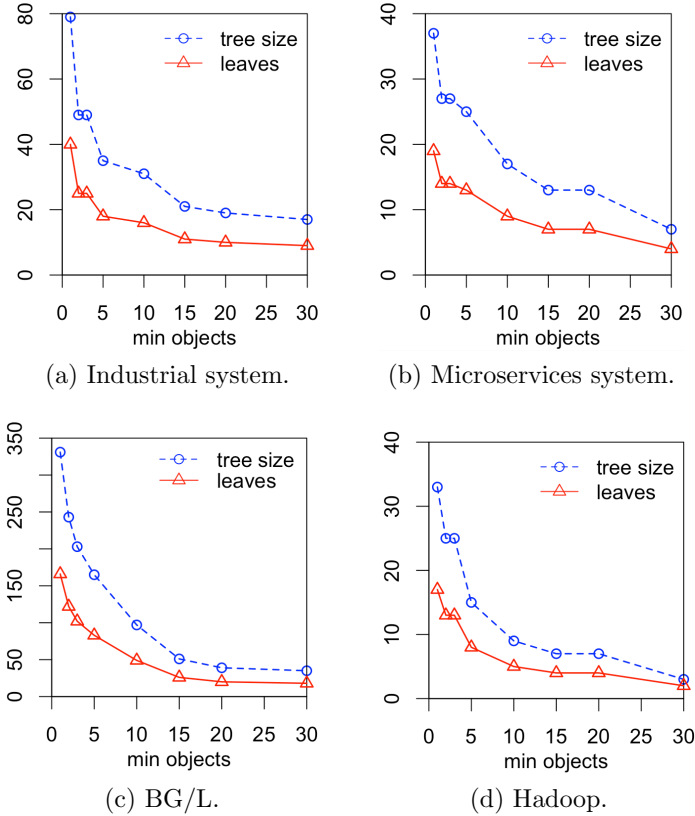


Figure 19: Sensitivity analysis of tree size and number of leaves of the decision tree with respect to the $minNumObj$ parameter.

the decision tree is a *supervised* classification technique. At training time it requires both normative and anomalous data points, and the availability of the labels in order to infer the decision predicates. This is a key limitation in applying decision trees to many real-life scenarios because training data will likely not cover all the potential anomalies that may occur in production. Analysis is done with the $J48$ tree implementation available with WEKA. For each dataset, we train and test the decision tree by means of a K fold cross validation approach. In consequence, the dataset is split into $K=10$ folds f_i ($1 \leq i \leq 10$), beforehand; for each fold f_i we train the decision tree with 9 folds f_j ($j \neq i$) and test it with the “held-out” fold f_i .

1028 **Hyperparameters selection.** The outcome of a decision tree may depend on the specific hyperparameters selected. Among the others, *minNumObj* is the hyperparameter that regulates the number of leaves that will be created when learning the tree: the lower the value, the higher the number
1032 of leaves and classification paths, thus causing overfitting. We perform a **sensitivity analysis** of the decision tree by varying *minNumObj*.

Figure 19 shows the sensitivity of tree-related figures –namely *tree size* and *number of leaves*– with respect to *minNumObj* for the systems. The parameter *minNumObj* is varied between 1 and 30; for each value of *minNumObj* we learn and test the decision tree by means of the *K fold* approach as explained above. We note that both *tree size* and *number of leaves* decrease as *minNumObj* increases. For example, in the industrial system dataset (Figure 19a) the number of leaves is 40 when *minNumObj*=1 and then it starts flattening at *minNumObj*=5; similarly, the number of leaves drops from 19 to 9 when *minNumObj* increases from 1 to 10 in the microservices system (Figure 19b). This issue is much more exacerbated in the BG/L dataset, which accounts for a larger number of attributes than the industrial and the microservices system. As shown in Figure 19c, at *minNumObj*=1 leaves and tree size are equal to 166 and 331; these figures drop sharply until they stabilize at 18 and 35, respectively, at *minNumObj*=30. As for Hadoop in Figure 19d, both tree size and number of leaves flatten at *minNumObj*=10. A large number of leaves indicates very *specialized* decision paths that tend to overfit the data and –in turn– to bias the evaluation metrics.

Figure 20 shows the evaluation metrics of the decision tree, i.e., F1 score, recall and precision, with respect to *minNumObj*. It can be noted that the metrics are strongly sensitive to changes of *minNumObj* for all the systems in hand. For example, in the industrial system the recall drops from 0.960 to 0.945 when *minNumObj* increases from 5 to 10 as shown in Figure 20a; similar considerations can be done for the microservices system and BG/L. Based on Figure 19, it can be reasonably stated that a “stable” tree model –where *tree size* and *number of leaves* stop changing sharply– is given by *minNumObj*=10. Whilst much of the work in the area overlooks this aspect, a stable tree provides a *non-overfitted* reflection of the data. Values of recall, precision and F1 score obtained for all the systems with the decision tree at *minNumObj*=10 are shown in Table 10. In all the cases, if not BG/L, recall, precision and F1 score are higher than 0.95; as for BG/L, the value of the metrics is 0.92.

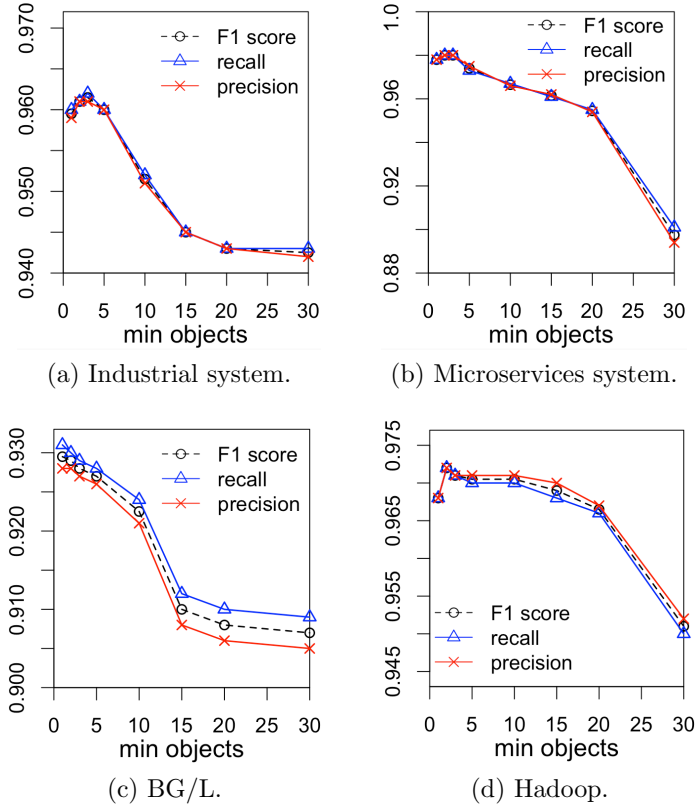


Figure 20: Sensitivity analysis of F1 score, recall, precision of the decision tree with respect to the $minNumObj$ parameter.

Table 10: Decision tree: evaluation metrics.

system	recall	precision	F1 score
industrial	0.95	0.95	0.95
microservices	0.96	0.96	0.96
BG/L	0.92	0.92	0.92
Hadoop	0.97	0.97	0.97

6.4. Vanilla AE

A *vanilla autoencoder* (AE) is an autoencoder that consists of only one hidden layer between the input and the output layer (Skansi, 2018). Typi-

1068

Table 11: Vanilla autoencoder: evaluation metrics.

system	recall	precision	F1 score
industrial	0.48	0.97	0.64
microservices	0.65	0.98	0.78
BG/L	0.81	0.91	0.86
Hadoop	0.98	0.95	0.97

Table 12: Variational autoencoder: evaluation metrics.

system	recall	precision	F1 score
industrial	0.71	0.97	0.82
microservices	0.76	0.94	0.84
BG/L	0.99	0.93	0.96
Hadoop	0.84	0.89	0.87

cally, the input layer has the same dimension as the output layer since the autoencoder attempts to reconstruct the input data point. Vanilla AEs have been used to address anomaly detection tasks. In our experiment we consider
1072 a vanilla AE with a hidden layer of 64 neurons; therefore, it includes $N-64-N$ neurons, where N is the number of logging entities for a given dataset. All other training hyperparameters are the same as AutoLog because they return the best result also for the vanilla AE. Table 11 shows the results obtained
1076 for the datasets. We observe that –applied to our datasets– the vanilla AE tends to achieve a high precision (e.g., 0.98 for the microservices system). On the other hand, recall is much lower; one notable exception is Hadoop where the recall is 0.98.

1080 6.5. Variational AE

A *variational autoencoder* (VAE) is similar to AEs as structure, but tries to obtain a better –continuous and not scattered– latent space by adopting a probabilistic distribution of each attribute in latent space (Kingma and
1084 Welling, 2019). The encoder stage of a VAE generates two vectors: a vector of means (μ) and a vector of variances (Σ); a vector z is then generated using the distribution $N(\mu_i, \Sigma_i)$, which will identify the points of the latent space. The addition of the *KullbackLeibler* (KL) divergence (van Erven and
1088 Harremos, 2014) to the cost function makes it possible to obtain a distribution as close as possible to a reference one, typically a Gaussian. The samples

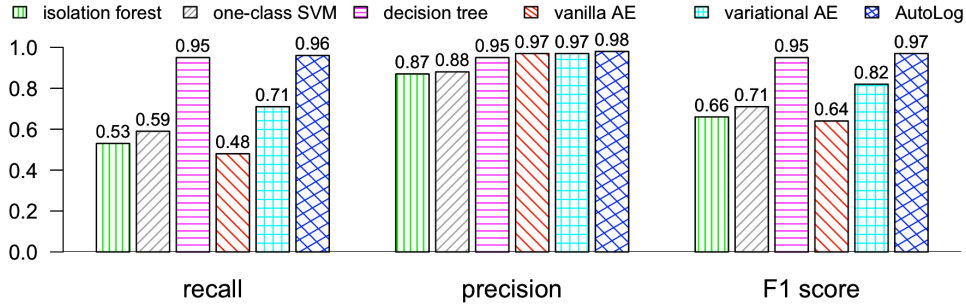


Figure 21: Industrial system.

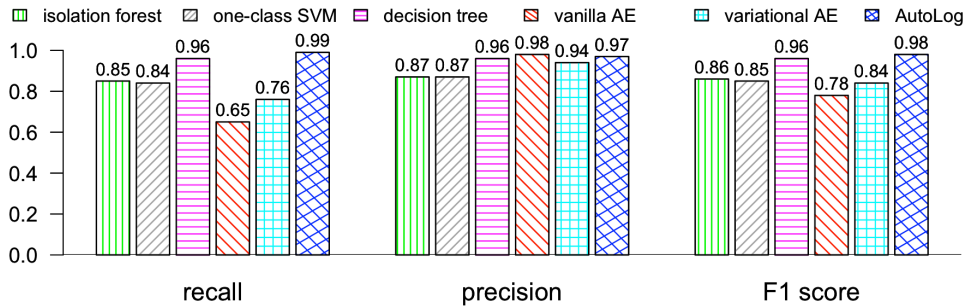


Figure 22: Microservices system.

obtained from this distribution are fed to a conventional decoder network.

For our VAE analysis, we select a latent layer of 2 neurons, while the encoder and the decoder layers have 256 neurons after tuning on the validation sets; as for the remaining training hyperparameters, they are the same of AutoLog because they returned the best results also for the VAE after tuning. Table 12 shows the results obtained for the reference systems. The best results are obtained with BG/L, where the VAE returns 0.99 recall and 0.93 precision. Overall, the precision of the VAE tends to be high across the systems –with the only exception of Hadoop– while recall is generally low.

6.6. Comparison of the Results

Barplots in Figure 21, 22, 23 and 24 show the evaluation metrics for all the techniques and systems addressed by our study. For each system, bars are grouped by evaluation metric in order to compare the different techniques; for each group of bars, the rightmost *cross-patterned* bar refers to the metrics of AutoLog presented in Table 6. Results reveal a mixture of findings, depending on the system and the technique. Isolation forest, one-class SVM

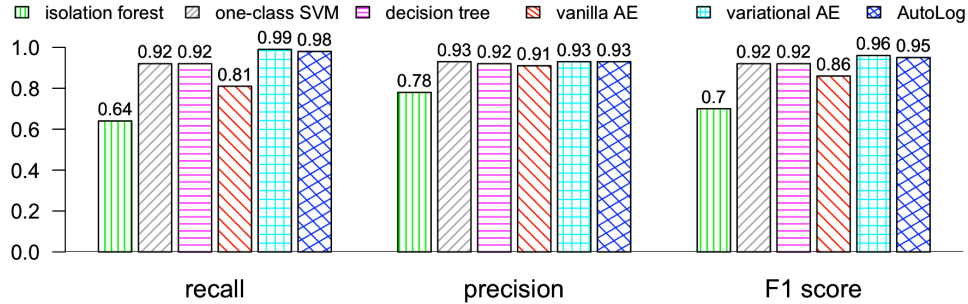


Figure 23: BG/L.

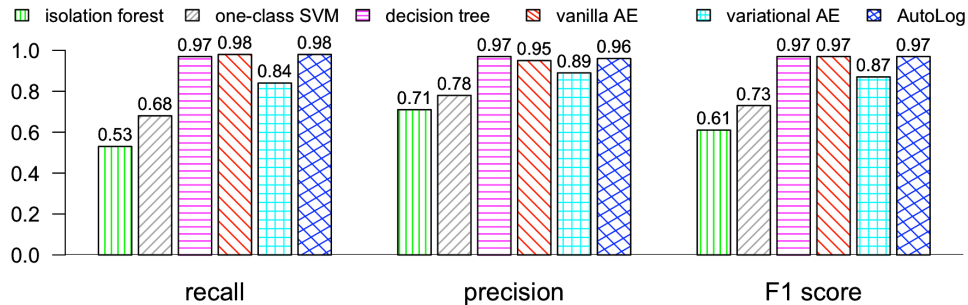


Figure 24: Hadoop.

and vanilla AE tend to achieve the lowest recall for all the systems; exceptions include (i) BG/L, where the recall of one-class SVM is 0.92 –thus equal
1108 to the decision tree– and (ii) Hadoop, where the vanilla AE performs as well as AutoLog. As said, the precision of vanilla AE is notably high for all the systems; however, its recall ranges between 0.48 (industrial system) and 0.98 (Hadoop), which makes it among the worst performing in terms of F1 score
1112 for all the systems if not Hadoop. As for the variational AE, it performs reasonably well only in the case of BG/L. AutoLog achieves among the highest values of the metrics for all the systems with few sporadic exceptions, as for example (i) the microservices systems, where the precision of AutoLog
1116 (0.97) is similar to the vanilla AE (0.98), or (ii) BG/L, where the F1 score of AutoLog (0.95) is similar to the variational AE (0.96). However, it should be noted that in spite of these sporadic exceptions, none of the techniques is able to fit all the systems when compared to AutoLog. Most notably,
1120 AutoLog is strongly competitive when compared to the decision tree, i.e., a typical supervised technique, that –different from AutoLog– requires both

normative and anomalous data points for training.

It is worth noting that, while the industrial and microservices systems were available at a private premise, BG/L and Hadoop are public datasets and widely-used benchmarks in log-based anomaly detection. In consequence, we can compare the metrics of AutoLog with other similar proposals in the literature. Authors in (Meng et al., 2019) propose a method –called LogAnomaly– and assess several state-of-the-art anomaly detection methods on the same BG/L dataset used in our paper, such as LogCluster and DeepLog. LogAnomaly is based on *template2Vec*, i.e., a template representation method to extract semantic and syntax information from log templates inspired by word embedding. According to the figures reported in (Meng et al., 2019) the F1 score of AutoLog in BG/L (0.95) is higher than DeepLog (0.93) and LogCluster (0.57) and similar to LogAnomaly (0.96). Unlike LogAnomaly, which achieves 0.94 recall and 0.97 precision, AutoLog achieves 0.98 recall and 0.93 precision; on the other hand, the precision of AutoLog is higher than DeepLog (0.90). As for the Hadoop dataset, the F1 score of AutoLog (0.97) is equal to the value achieved on the same dataset by the recent OneLog approach (Hashemi and Mäntylä, 2021). Overall, AutoLog is inline with other anomaly detection techniques assessed with the reference BG/L and Hadoop datasets.

Coming to the comparison with approaches implying sequential embedding, the work (Guo et al., 2021) applies BERT to the BG/L log and achieves 0.92 recall and 0.89 precision: both the figures are much lower than AutoLog applied to the same dataset. As for other related approaches using BERT with BG/L, it is worth mentioning (Hirakawa et al., 2020) and (Li et al., 2020). The former uses BERT-Base as a transformer encoder and its F1 score on BG/L is 0.89, thus lower than AutoLog: more important, a conventional LSTM autoencoder was demonstrated to achieve a much lower F1 score equals to 0.85 (Hirakawa et al., 2020). The latter proposes SwissLog (Li et al., 2020), which is applied in conjunction with LSTM, Bidirectional LSTM (BiLSTM) and Attention-based BiLSTM. The F1 score of these three variants on BG/L is 0.95, 0.96 and 0.99, respectively: the F1 score of AutoLog is in line with the first two variants, although lower than the Attention-based BiLSTM. Nevertheless, it must be noted that SwissLog is conceived and tested with a much more narrow fault model, consisting of those faults resulting in log sequence order changes and log time interval changes: different from SwissLog, AutoLog is not constrained by specific types of faults.

7. Limitations and Threats to Validity

1160 As for many existing anomaly detection techniques, AutoLog relies on
the availability of normative logs; moreover, AutoLog does not need to know
anomalies at training time. Whilst our approach is feasible in practice, we
are aware that **assembling normative logs** is a complex matter and may
1164 depend on the specific system in hand. For example, both the industrial and
microservices systems were deployed in a dedicated LAN environment along
with client applications supplied by the vendor and representative bench-
marks, which makes us confident that logs actually reflected normative con-
1168 ditions. In practice, normative logs might accidentally account for anomalous
events, which are hard to be filtered out before training. In this respect,
the 90th percentile approach –used to determine the detection threshold of
AutoLog– aims to mitigate the impact of measurement outliers due to ac-
1172 cidental anomalies in the normative logs. It is worth noting that in BG/L
we deal with spontaneous anomalies with much less confidence on the genu-
ineness of the normative logs. In this respect, the F1 score obtained by
AutoLog for BG/L is lower than the industrial and microservices systems,
1176 although inline with other anomaly detection techniques assessed with the
same dataset.

At the current stage of deployment, AutoLog is a “partial” fit for the
detection of **slow attacks**, whose activity may span multiple chunks. In fact,
1180 while AutoLog can detect the activity occurring within individual chunks, it
misses the ability to “connect the dots” across multiple chunks collected at
different times. A slow attack may cause AutoLog to generate multiple alerts,
which are supposed to be reconciled later on. On the other hand, AutoLog
1184 can detect the activity affecting multiple chunks collected at the same time.

Regarding the **anomaly detection** technique implemented by AutoLog,
semi-supervised training tends to obtain from the AE high-quality recon-
structions of normative inputs, in such a way that anomalies can be identi-
1188 fied by a higher reconstruction error. However, recent studies point out that
sometimes autoencoders can provide good reconstructions also for outliers,
which are misclassified as being in-distribution. This is obviously detrimen-
tal for anomaly detection. Possible solutions tend to improve the anomaly
1192 detection capabilities of a simple autoencoder design by resorting to *autoen-
coder ensembles* (randomly connected autoencoders with different structures
and connection densities (Chen et al., 2017)), discriminator networks with a
direct anomaly score (Tong et al., 2020), or the use of the Grubbs and the

1196 PauTa criterion to identify the reconstruction errors corresponding to the
outliers based on the traditional threshold method (Wan et al., 2019). In
our experiments, which involved the use of very heterogeneous data sources,
the problem did not arise –or was very limited in size– as documented by
1200 the good detection performance obtained though our AE. However, possible
outlier misclassification is an issue to be considered for our future work to
optimize the current design.

As for any data-driven study, there may be concerns regarding the validity
1204 and generalizability of the results. We discuss them based on the four aspects
of validity listed in (Wohlin et al., 2000).

Construct validity. The study builds around the intuition that numeric
scores computed from system logs can be used for detecting anomalies of com-
1208 puter systems. This *construct* has been investigated in the context of four
systems: an industrial system in the transportation domain, a microservices-
based installation implementing a standard multimedia architecture adopted
by large telcos, up to publicly-available system logs –common benchmarks
1212 in the literature– from a BG/L supercomputer and a Hadoop cluster. The
industrial and microservices systems are run with representative load gen-
erators. The BG/L log is a public dataset and accounts for hardware and
software errors observed over 215 days of operations at LLNL; Hadoop data
1216 contain different types of service failures. The study is supported by ex-
tensive experimentation leveraging widely-consolidated statistical methods,
deep learning framework and evaluation metrics.

Internal validity. The results and key findings of this paper are based
1220 on direct measurement experiments, where we analyze systems logs from the
reference systems obtained under a mixture of simulated and spontaneous
normative operations and anomalies. For example, simulated anomalies in
the industrial and microservices systems consist of bruteforce authentica-
1224 tion attempts, tampering and misuse, which are inspired by widely-accepted
taxonomies in the area. We rely on well-founded log analysis methods for
extracting quantitative scores from logs. Noteworthy, AutoLog –based solely
on the knowledge of the normative data points– is not biased by the anoma-
1228 lies in hand. The use of such diverse system logs and anomalies aims to
mitigate internal validity threats.

Conclusion validity. Conclusions have been inferred by assessing four
independent system logs and the sensitivity of key results with respect to the
1232 experimental choices. For example, we analyze the impact of different au-
toencoder configurations on the reconstruction error; similarly, the decision

tree is assessed by varying a critical model parameter, such as *minNumObj*. Experiments are complemented by a preliminary discussion on PCA and clustering to gain insights into the challenges existing in our domain. More importantly, we compare AutoLog with a wide set of techniques including isolation forest, one-class SVM, decision trees, vanilla autoencoder and variational autoencoder. We present an extensive discussion of the results. The key findings of the study are consistent across the datasets, which provides a reasonable level of confidence on the analysis.

External validity. The steps of our analysis can be applied to other systems. Nowadays, logs are ubiquitously emitted by almost any system and there exists a wide range of tools for storing and handling logs, which make our approach definitively feasible in practice. In fact, in this paper we successfully ported the experiments across four independent systems –emitting heterogeneous logs– to mitigate external validity threats. Our analysis approach does not interfere with system operations: as only information from logs is used, the approach is inherently *non intrusive*. We are confident that the experimental details provided in the paper would support the replication of our study by future researchers and practitioners.

8. Conclusion and Future Work

This paper proposed AutoLog, a novel approach for anomaly detection based on deep autoencoding of system logs. AutoLog aims to overcome the challenges of existing log management technologies for the analysis of built-in and proprietary logs files, which lack standard formats. More importantly, AutoLog capitalizes on semi-supervised learning, which does not require anomalies during training. This is potentially valuable to complement current technologies that rely on pre-established specifications of anomalies. We have conducted an extensive experimentation in the context of four systems and compared AutoLog with a wide set of techniques. The recall of AutoLog ranged between 0.96 and 0.99, while its precision was within 0.93 and 0.98, depending on the system. AutoLog is strongly competitive if compared to a typical supervised technique, such as decision trees; as for BG/L and Hadoop –used as benchmarks by related papers– AutoLog is inline with other log-based anomaly detection techniques available in the literature.

In the future, we will extend our analysis to further systems as well as to existing security datasets, in order to understand the limitations of AutoLog and potential mitigations. At the time being, there are several **open**

challenges in anomaly detection that reflect on AutoLog. For example, the availability and construction of normative baselines is a complex matter in log analysis, given the variability of real-life workload and systems. Moreover, log analysis is a “moving” target: software upgrades or changes to the configurations may alter meaning and character of the logs. As for the detection technique used by AutoLog, one key challenge is the threshold selection technique and how it is affected by accidental anomalies intertwined with the normative data. In this respect, future research will investigate some countermeasures to mitigate these issues as well as other threshold selection methods. Other important research avenues pertain to potential strategies to partition logs into chunks, such as considering overlapping chunks, and the explicative power of the autoencoder. Future research will also investigate the actions attackers might take to evade detection and the detection of slow attacks. From a technical standpoint, we will address the implementation of accessory components for AutoLog, such a re-training mechanism, and the use of plugins to extend traditional SIEM’s capabilities with the proposed anomaly detection method.

References

- Adiga, N. R. et al. (2002). An overview of the BlueGene/L supercomputer. In *Proc. Conference on Supercomputing*, pages 1–22. IEEE.
- Almotiri, J., Elleithy, K., and Elleithy, A. (2017). Comparison of autoencoder and principal component analysis followed by neural network for e-learning using handwritten recognition. In *Proc. Long Island Systems, Applications and Technology Conference*, pages 1–5. IEEE.
- Aygun, R. C. and Yavuz, A. G. (2017). Network Anomaly Detection with Stochastically Improved Autoencoder Based Models. In *Proc. International Conference on Cyber Security and Cloud Computing*, pages 193–198. IEEE.
- Bertero, C., Roy, M., Sauvanaud, C., and Tredan, G. (2017). Experience report: Log mining using natural language processing and application to anomaly detection. In *Proc. International Symposium on Software Reliability Engineering*, pages 351–360. IEEE.

- 1304 Bhatt, S., Manadhata, P. K., and Zomlot, L. (2014). The operational role of security information and event management systems. *IEEE Security Privacy*, 12(5):35–41.
- Campos, J. R., Vieira, M., and Costa, E. (2018). Exploratory study of machine learning techniques for supporting failure prediction. In *Proc. European Dependable Computing Conference*, pages 9–16. IEEE.
- 1308 Carrington, A. M., Manuel, D. G., Fieguth, P. W., Ramsay, T., Osmani, V., Wernly, B., Bennett, C., Hawken, S., McInnes, M., Magwood, O., Sheikh, Y., and Holzinger, A. (2021). Deep ROC analysis and AUC as balanced average accuracy to improve model selection, understanding and interpretation. arXiv:2103.11357.
- 1312 Catillo, M., Rak, M., and Villano, U. (2020). 2L-ZED-IDS: A two-level anomaly detector for multiple attack classes. In *Proc. Web, Artificial Intelligence and Network Applications*, Advances in Intelligent Systems and Computing, pages 687–696. Springer.
- Chen, J., Sathe, S., Aggarwal, C., and Turaga, D. (2017). Outlier detection with autoencoder ensembles. In *Proc. SIAM International Conference on Data Mining*, pages 90–98. SIAM.
- 1320 Cinque, M., Cotroneo, D., Della Corte, R., and Pecchia, A. (2016). Characterizing direct monitoring techniques in software systems. *IEEE Transactions on Reliability*, 65(4):1665–1681.
- 1324 Cinque, M., Cotroneo, D., and Pecchia, A. (2018). Challenges and directions in security information and event management (SIEM). In *Proc. International Symposium on Software Reliability Engineering Workshops*, pages 95–99. IEEE.
- 1328 Du, J., Vong, C., Pun, C., Wong, P., and Ip, W. (2017a). Post-boosting of classification boundary for imbalanced data using geometric mean. *Neural Networks*, 96:101–114.
- 1332 Du, M., Li, F., Zheng, G., and Srikumar, V. (2017b). DeepLog: Anomaly detection and diagnosis from system logs through deep learning. In *Proc. Conference on Computer and Communications Security*, pages 1285–1298. ACM.

- 1336 Fahimeh, F. and Heikkonen, J. (2018). A Deep Auto-Encoder based Approach for Intrusion Detection System. In *Proc. International Conference on Advanced Communications Technology*, pages 178–183. IEEE.
- Farshchi, M., Schneider, J.-G., Weber, I., and Grundy, J. (2018). Metric selection and anomaly detection for cloud operations using log and metric correlation analysis. *Journal of Systems and Software*, 137:531 – 549.
- 1340 Farzad, A. and Gulliver, T. A. (2020). Unsupervised log message anomaly detection. *ICT Express*, 6(3):229–237.
- Farzad, A. and Gulliver, T. A. (2021). Log Message Anomaly Detection and Classification Using Auto-B/LSTM and Auto-GRU. arXiv:1911.08744.
- 1344 Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep Learning*. MIT Press.
- Guo, H., Yuan, S., and Wu, X. (2021). LogBERT: Log Anomaly Detection via BERT. arXiv:2103.04475.
- 1348 Hansen, J. and Siewiorek, D. (1992). Models for time coalescence in event logs. In *Proc. International Symposium on Fault-Tolerant Computing*, pages 221–227. IEEE.
- 1352 Hashemi, S. and Mäntylä, M. (2021). OneLog: Towards end-to-end training in software log anomaly detection. arXiv:2104.07324.
- Hawkins, S., He, H., Williams, G., and Baxter, R. (2002). Outlier Detection Using Replicator Neural Networks. In *Proc. International Conference on Data Warehousing and Knowledge Discovery*, pages 170–180. Springer.
- 1356 He, P., Zhu, J., He, S., Li, J., and Lyu, M. R. (2016a). An evaluation study on log parsing and its use in log mining. In *Proc. International Conference on Dependable Systems and Networks*, pages 654–661. IEEE.
- 1360 He, S., Zhu, J., He, P., and Lyu, M. R. (2016b). Experience Report: System Log Analysis for Anomaly Detection. In *Proc. International Symposium on Software Reliability Engineering*, pages 207–218. IEEE.
- Hinton, G. E., Osindero, S., and Teh, Y.-W. (2006). A Fast Learning Algorithm for Deep Belief Nets. *Neural Computation*, 18(7):1527–1554.

- 1364 Hirakawa, R., Tominaga, K., and Nakatoh, Y. (2020). Software log anomaly detection through one class clustering of transformer encoder representation. In Stephanidis, C. and Antona, M., editors, *Proc. HCI International - Posters*, pages 655–661. Springer.
- 1368 Holzinger, A., Hörtenhuber, M., Mayer, C., Bachler, M., Wassertheurer, S., Pinho, A. J., and Koslicki, D. (2014). On entropy-based data mining. In Holzinger, A. and Jurisica, I., editors, *Interactive Knowledge Discovery and Data Mining in Biomedical Informatics: State-of-the-Art and Future Challenges*, pages 209–226. Springer.
- 1372 Kim, S., Jo, W., and Shon, T. (2020). APAD: Autoencoder-based payload anomaly detection for industrial IoE. *Applied Soft Computing*, 88:106017.
- Kingma, D. P. and Welling, M. (2019). An introduction to variational autoencoders. *Foundations and Trends in Machine Learning*, 12(4):307–392.
- 1376 Li, X., Chen, P., Jing, L., He, Z., and Yu, G. (2020). Swisslog: Robust and unified deep learning based log anomaly detection for diverse faults. In *Proc. International Symposium on Software Reliability Engineering*, pages 92–103. IEEE.
- 1380 Lin, Q., Zhang, H., Lou, J.-G., Zhang, Y., and Chen, X. (2016). Log clustering based problem identification for online service systems. In *Proc. International Conference on Software Engineering Companion*, pages 102–111. IEEE.
- 1384 Liu, D., Zhao, Y., Xu, H., Sun, Y., Pei, D., Luo, J., Jing, X., and Feng, M. (2015). Opprentice: Towards practical and automatic anomaly detection through machine learning. In *Proc. Internet Measurement Conference*, pages 211–224. ACM.
- 1388 Liu, F. T., Ting, K. M., and Zhou, Z. (2008). Isolation forest. In *Proc. International Conference on Data Mining*, pages 413–422. IEEE.
- 1392 Liu, J., Song, K., Feng, M., Yan, Y., Tu, Z., and Zhu, L. (2021). Semi-supervised anomaly detection with dual prototypes autoencoder for industrial surface inspection. *Optics and Lasers in Engineering*, 136:106324.

- 1396 Lu, S., Rao, B., Wei, X., Tak, B., Wang, L., and Wang, L. (2017). Log-based abnormal task detection and root cause analysis for spark. In *Proc. International Conference on Web Services*, pages 389–396. IEEE.
- Maci-Fernndez, G., Camacho, J., Magn-Carrin, R., Garca-Teodoro, P., and Thern, R. (2018). UGR’16: A new dataset for the evaluation of cyclostationarity-based network IDSs. *Computers & Security*, 73:411–424.
- 1400 Meng, W., Liu, Y., Zhu, Y., Zhang, S., Pei, D., Liu, Y., Chen, Y., Zhang, R., Tao, S., Sun, P., and Zhou, R. (2019). LogAnomaly: Unsupervised Detection of Sequential and Quantitative Anomalies in Unstructured Logs. In *Proc. International Joint Conf. on Artificial Intelligence*, pages 4739–4745.
- 1404 International Joint Conferences on Artificial Intelligence Organization.
- Miller, D., Harris, S., Harper, A., VanDyke, S., and Blask, C. (2010). *Security Information and Event Management (SIEM) Implementation*. McGraw-Hill Education.
- 1408 Nguyen, Q. P., Lim, K. W., Divakaran, D. M., Low, K. H., and Chan, M. C. (2019). GEE: A gradient-based explainable variational autoencoder for network anomaly detection. In *Proc. Conference on Communications and Network Security*, pages 91–99. IEEE.
- 1412 Oliner, A., Ganapathi, A., and Xu, W. (2012). Advances and challenges in log analysis. *Communications of the ACM*, 55(2):55–61.
- Oliner, A. and Stearley, J. (2007). What Supercomputers Say: A Study of Five System Logs. In *Proc. International Conference on Dependable Systems and Networks*, pages 575–584. IEEE.
- 1416
- Oprea, A., Li, Z., Yen, T., Chin, S. H., and Alrwais, S. (2015). Detection of early-stage enterprise infection by mining large-scale log data. In *Proc. International Conference on Dependable Systems and Networks*, pages 45–56. IEEE.
- 1420
- Pang, G., Shen, C., Cao, L., and Hengel, A. V. D. (2021). Deep Learning for Anomaly Detection: A Review. *ACM Computing Surveys*, 54(2):1–38.
- Pisner, D. A. and Schnyer, D. M. (2020). Chapter 6 - Support vector machine.
- 1424 In *Machine Learning*, pages 101–121. Academic Press.

- Qian, Y., Ying, S., and Wang, B. (2020). Anomaly detection in distributed systems via variational autoencoders. In *Proc. International Conference on Systems, Man, and Cybernetics*, pages 2822–2829. IEEE.
- 1428 Quan, X., Wenyin, L., and Qiu, B. (2011). Term weighting schemes for question categorization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(5):1009–1021.
- Ruff, L., Vandermeulen, R. A., Grnitz, N., Binder, A., Mller, E., Mller, K.-R., and Kloft, M. (2020). Deep semi-supervised anomaly detection. In *Proc. International Conference on Learning Representations*.
- Ruii, D. (1999). *Cautionary tales: stealth coordinated attack how to*. <http://www.ouah.org/stealthhowto.html>.
- 1436 Sakurada, M. and Yairi, T. (2014). Anomaly detection using autoencoders with nonlinear dimensionality reduction. In *Proc. Workshop on Machine Learning for Sensory Data Analysis*, pages 4–11. ACM.
- Salton, G. and Buckley, C. (1988). Term-weighting approaches in automatic text retrieval. *Information Processing & Management*, 24(5):513 – 523.
- 1440 Schölkopf, B., Platt, J. C., Shawe-Taylor, J. C., Smola, A. J., and Williamson, R. C. (2001). Estimating the support of a high-dimensional distribution. *Neural Computing*, 13(7):14431471.
- 1444 Shone, N., Ngoc, T. N., Phai, V. D., and Shi, Q. (2018). A Deep Learning Approach to Network Intrusion Detection. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2(1):41–50.
- Skansi, S. (2018). *Autoencoders*, pages 153–163. Springer International Publishing.
- 1448 Song, C., Liu, F., Huang, Y., Wang, L., and Tan, T. (2013). Auto-encoder based data clustering. In *Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications*, pages 117–124. Springer.
- 1452 Stearley, J. and Oliner, A. J. (2008). Bad words: Finding faults in Spirit’s syslog. In *Proc. International Symposium on Cluster Computing and the Grid*, pages 765–770. IEEE.

- 1456 Su, Y., Zhao, Y., Niu, C., Liu, R., Sun, W., and Pei, D. (2019). Robust Anomaly Detection for Multivariate Time Series through Stochastic Recurrent Neural Network. In *Proc. International Conference on Knowledge Discovery & Data Mining*, pages 2828–2837. ACM.
- 1460 Svacina, J., Raffety, J., Woodahl, C., Stone, B., Cerny, T., Bures, M., Shin, D., Frajtak, K., and Tisnovsky, P. (2020). On Vulnerability and Security Log analysis: A Systematic Literature Review on Recent Trends. In *Proc. International Conference on Research in Adaptive and Convergent Systems*, pages 175–180. ACM.
- 1464 Tong, A., Wolf, G., and Krishnaswamy, S. (2020). Fixing bias in reconstruction-based anomaly detection with Lipschitz discriminators. In *Proc. International Workshop on Machine Learning for Signal Processing*, pages 1–6. IEEE.
- 1468 van Erven, T. and Harremoës, P. (2014). Rényi divergence and Kullback-Leibler divergence. *IEEE Transactions on Information Theory*, 60(7):3797–3820.
- 1472 Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., and Manzagol, P. A. (2010). Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11:3371–3408.
- 1476 Wadekar, A., Gupta, T., Vijan, R., and Kazi, F. (2019). Hybrid CAE-VAE for Unsupervised Anomaly Detection in Log File Systems. In *Proc. International Conference on Computing, Communication and Networking Technologies*, pages 1–7. IEEE.
- 1480 Wan, F., Guo, G., Zhang, C., Guo, Q., and Liu, J. (2019). Outlier Detection for Monitoring Data Using Stacked Autoencoder. *IEEE Access*, 7:173827–173837.
- 1484 Wohlin, C., Runeson, P., Höst, M., Ohlsson, M. C., Regnell, B., and Wesslén, A. (2000). *Experimentation in Software Engineering: An Introduction*. Kluwer Academic.
- Xu, W., Huang, L., Fox, A., Patterson, D. A., and Jordan, M. I. (2008). Mining console logs for large-scale system problem detection. In *Proc.*

- 1488 *Tackling Computer Systems Problems with Machine Learning Techniques*, pages 4–4. USENIX.
- Yadav, R. B., Kumar, P. S., and Dhavale, S. V. (2020). A Survey on Log Anomaly Detection using Deep Learning. In *Proc. International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions)*, pages 1215–1220. IEEE.
- 1492 Yang, R., Qu, D., Gao, Y., Qian, Y., and Tang, Y. (2019). nLSALog: An anomaly detection framework for log sequence in security management. *IEEE Access*, 7:181152–181164.
- 1496 Yuan, Y., Srikant Adhatarao, S., Lin, M., Yuan, Y., Liu, Z., and Fu, X. (2020). ADA: Adaptive deep log anomaly detector. In *Proc. INFOCOM - Conference on Computer Communications*, pages 2449–2458. IEEE.
- 1500 Zhang, K., Xu, J., Min, M. R., Jiang, G., Pelechrinis, K., and Zhang, H. (2016). Automated IT system failure prediction: A deep learning approach. In *Proc. International Conference on Big Data*, pages 1291–1300. IEEE.
- Zhang, Y., Lu, Z., and Wang, S. (2021). Unsupervised feature selection via transformed auto-encoder. *Knowledge-Based Systems*, 215:106748.
- 1504 Zhao, Y., Hao, K., Tang, X., Chen, L., and Wei, B. (2021). A conditional variational autoencoder based self-transferred algorithm for imbalanced classification. *Knowledge-Based Systems*, 218:106756.
- 1508 Zoppi, T., Ceccarelli, A., and Bondavalli, A. (2016). Context-awareness to improve anomaly detection in dynamic service oriented architectures. In Skavhaug, A., Guiochet, J., and Bitsch, F., editors, *Computer Safety, Reliability, and Security*, pages 145–158. Springer.