Intrusion Detection and Datasets

Subjects: Computer Science, Artificial Intelligence

Contributor: Joaquín Gaspar Medina-Arco , Roberto Magán-Carrión , Rafael Alejandro Rodríguez-Gómez , Pedro García-Teodoro

With the significant increase in cyber-attacks and attempts to gain unauthorised access to systems and information, Network Intrusion-Detection Systems (NIDSs) have become essential detection tools. Anomaly-based systems use machine learning techniques to distinguish between normal and anomalous traffic. They do this by using training datasets that have been previously gathered and labelled, allowing them to learn to detect anomalies in future data. However, such datasets can be accidentally or deliberately contaminated, compromising the performance of NIDS.

anomaly detection

NIDS

deep learning

datasets

network traffic labelling

1. Introduction

Network Intrusion-Detection Systems (NIDSs) represent a primary cybersecurity mechanism for identifying potential attacks on a communication network. To accomplish this goal, they analyse the network traffic passing through the system, regardless of whether it is internally generated or originated from external entities targeting the network. Detecting intrusions allows network administrators to become aware of system vulnerabilities and to make quick decisions to abort or mitigate attacks. Additionally, NIDSs allow them to implement measures to strengthen the system in the future ^[1].

NIDSs can be categorised into various typologies based on two fundamental principles: architecture and techniques employed. Focusing on the architecture, NIDS can be classified as host-based, network-based, and collaborative approaches between different components. According to the detection technique, the classification may be signature-based, Stateful Protocol Analysis-based, or anomaly detection-based NIDSs ^[2].

Signature-based NIDSs possess a repository of network patterns representing prevalent network attacks. Their operating mode is to match the network sequence they examine with their knowledge base to detect potential attacks ^[3].

Alternatively, Stateful Protocol Analysis-based NIDSs rely on their comprehensive understanding of the monitored protocol. They analyse all interactions to identify a sequence of actions that might result in a vulnerability or insecurity ^[3].

In contrast, anomaly-detection-based NIDSs employ mechanisms to detect abnormal network traffic behaviour. These anomalous activities typically correspond to network traffic patterns that have a significantly low likelihood of occurring or are markedly misaligned with normal traffic. Acutely objective, anomaly detection allows for the handling of novel or previously unknown attacks (*zero days*). This is because such attacks generate traffic patterns that have not been found before, and this type of NIDS often relies on the use of machine learning techniques to carry out anomaly detection. When this approach is followed, the subjective evaluation of attacks is effectively circumvented.

Different strategies have been employed to detect anomalies in NIDS through various machine learning techniques ^{[4][5]}, including statistical techniques like Principal Component Analysis (PCA) ^[6] or Markov models ^{[7][8]}; classification techniques like Artificial Neural Networks (ANNs) ^{[9][10][11][12]}, Support Vector Machines (SVMs) ^[6], deep learning models ^{[13][14]} including Autoencoders ^{[9][15]}, or Decision Trees including Random Forest ^[16]; and clustering like outlier detection ^[17]. Using these techniques requires a multi-perspective approach to tackling the problem, which can be categorised as supervised, semi-supervised, or unsupervised, depending on the specific technique chosen ^[18].

Regardless of the technique used for anomaly detection in NIDS, the underlying models must be trained to distinguish normal traffic from anomalous traffic. This training process utilises datasets comprising real, synthetic, or a combination of both network traffic. To be more concise,

- **Synthetic traffic datasets** are created by generating traffic in a controlled environment that emulates a realworld setting. The generated traffic may include traffic related to known attacks, providing enough samples for machine learning models to competently identify and detect such anomalies. This enables the optimisation of the dataset regarding the size and balance between regular and irregular traffic samples. It also ensures the correct labelling of each observation as it has been intentionally and deliberately generated. Such observations can be, for instance, the traffic flows seen in the network. However, a potential issue is that it may not accurately reflect the network traffic patterns observed in a genuine environment.
- Real traffic datasets capture all network communications within a real productive environment. This implies access to the patterns of network traffic consumption and usage that take place in an actual scenario and potentially any cyber-attacks that may occur. Unlike synthetic datasets, real traffic samples may be biased or imbalanced, with the presence of anomalous traffic often being minimal or completely absent. It is necessary to carry out a subsequent process to assign a normality or attack label to each flow for its use in machine learning models during training phases.
- **Composite datasets** are the ones generated by combining real environment data and synthetic traffic to introduce attack patterns.

Regardless of the AI model used in a NIDS, the dataset's labelling accuracy is crucial to maintaining high model performance. This principle applies equally to supervised and unsupervised learning. In supervised learning,

labelling is necessary to enable models to learn how to identify anomalous traffic. In contrast, unsupervised learning generally assumes that the training dataset consists of normal traffic only and is, therefore, free of anomalies.

2. Datasets for Network Security Purposes

To effectively train any AI model, especially those constituting NIDSs based on anomaly detection, a prerequisite is a comprehensive dataset. This dataset should encompass a sufficient number of samples that represent all the various classes or patterns, whether benign or malicious. This foundational dataset enables the model to learn and predict accurately during subsequent training phases. In the specific case of NIDSs, a large and correctly labelled dataset is assumed ^[19]. The quality of the trained models depends to some extent on the quality of the data on which they were trained ^[20], so it is important to make a thorough analysis of the typology of datasets available in the NIDS domain.

Before reviewing the different datasets available in the field of cybersecurity, it is necessary to define the criteria according to which these datasets will be analysed:

- Availability: Understood as free access (Public) to the dataset or, on the contrary, of reserved access, by means of payment or explicit request (Protected).
- **Collected data**: Some datasets collect traffic packet for each packet (e.g., PCAP files), others collect information associated with traffic flows between devices (e.g., NetFlow), and others extract features from the flows by combining them with data extracted from the packets.
- Labelling: This refers to whether each observation in the dataset has been identified as normal, anomalous, or even belonging to a known attack. Or, conversely, no labelling is available, in which case they are intended for unsupervised learning models.
- **Type**: The nature of a dataset may be synthetic, where the process and environment in which the dataset is generated are controlled, or it may be the result of capturing traffic in a real environment.
- **Duration**: Network traffic datasets consist of network traffic recorded over a specific time interval, which may range from hours to days, months, or even years.
- Size: the depth of the dataset in terms of the number of records or the physical size and their distribution across the different classes.
- Freshness: It is also important to consider the year in which the dataset was created, as the evolution of attacks and network usage patterns may not be reflected in older datasets, thus compromising their validity in addressing current issues.

A summary of the datasets analysed according to the characteristics described above is shown in **Table 1**.

Dataset	Availability	Collected Data	Labeled	Туре	Duration *	Size **	Year	Freshness B	alanced
DARPA ^[21]	Public	packets	yes	synthetic	7 weeks	6.5TB	1998— 1999	questioned	no
NSL-KDD [<u>22</u>]	Public	features	yes	synthetic	N.S.	5M o.	1998— 1999	questioned	yes
Kyoto 2006+ ^[23]	Public	features	yes	real	9 years	93M 0.	2006– 2015	yes	yes
Botnet ^[24]	Public	packets	yes	synthetic	N.S.	14GB р.	2010– 2014	yes	yes
UNSW- NB15 ^[25]	Public	features	yes	synthetic	31 hours	2.5M 0.	2015	yes	no
UGR'16 ^[26]	Public	flows	yes	real	6 months	17B f.	2016	yes	no
CICIDS2017 [27]	Protected	flows	yes	synthetic	5 days	3.1M f.	2017	yes	no
IDS2018 ^[28]	Protected	features	yes	synthetic	10 days	1M o.	2018	yes	no
NF-UQ- NIDS ^[29]	Public	flows	yes	synthetic	N.S.	12M f.	2021	yes	no

Table 1. Overview of available network datasets.

References

* N.S. means not specified. ** Expressed in flows (f.), observations (o.), or packets (p.). An observation denotes a 1. Ahmad, Z.; Shahid Khan, A.; Waia Shiang with; All solution detection

system: A systematic study of machine learning and deep learning approaches. Trans. Emerg. **2.1 EDARPA Datasets** 2021, 32, e4150.

Ereliae, by MitrBichardh-Liab Gattory, the BARPA-Uga Kets, With Riod datasetis, any stemars from host where yes and in the view of interview of interview of interview of anomaly detection to the training subset consists of only three weeks of observations. In both cases, two weeks of observed network traffic is reserved for validation. All observations are 4. Patcha, A.; Park, J.M. An overview of anomaly detection techniques: Existing solutions and latest labeled and contain a total of 200 observations of up to 58 attacks of different typologies, including the typologies of the technologies of the typologies, port scanning, and user-to-root (U2R) or remote-to-local attacks (R2L) [21].

Breserrates doe opife; the azet of the power which are emilanded to any via zavious scanaous about a set works seeints usion poletect 12 h a Teologi ques, avesterns sandes tallequessio Correip telli Steicu 122009, 28, 18-28.

6. Wang, H.; Gu, J.; Wang, S. An effective intrusion detection framework based on SVM with feature **2.2. KDD Dataset** augmentation. Knowl.-Based Syst. 2017, 136, 130–139.

KDD22 [22] iDaydataset, created fobthsethirthutsunational Knowledge Disnerance and Data Minipaga Toolsa Production based then the DARPAt dataset. 34 like gho latter, KDD99 is a dataset whose format is based on the extraction of features (up to 41 [33]) from network flows rather than the recording of raw observed data. It is a synthetic dataset 8. Mahoney, M.V.: Chan, P.K. Learning nonstationary models of normal network traffic for detecting but takes into account the actual traffic observed in military network environments. Access to the dataset is open, novel attacks. In Proceedings of the Eighth ACM SIGKDD International Conference on and, despite its longevity, it is still available. In terms of size, the dataset contains almost 5 million observations, Knowledge Discovery and Data Mining, Edmonton, AB, Canada, 23–26 July 2002; pp. 376–385. including the same typology of attacks as DARPA, i.e., DoS, port scanning and privilege escalation attacks.

9. Mirsky, Y.; Doitshman, T.; Elovici, Y.; Shabtai, A. Kitsune: An Ensemble of Autoencoders for Online Similar to RAR PAR a line of the solution a wide zone play of dataset of the second se Specifically, concerns have been raised about the lack of consistency between the number of attack types in the 10. Li, J.: Manikopoulos, C.: Jorgenson, J.: Ucles, J. HUDE: A Hierarchical Network Intrusion Detection training subset and those available in the validation subset. Additionally, the dataset is deemed outdated in the System Using Statistical Preprocessing and Neural Network Classification. In Proceedings of the context of contemporary world communications.

2001 IEEE Workshop on Information Assurance and Security, West Point, NY, USA, 5–6 June

2.3.00SL-KDD Dataset

Poojitha, G.; Kumar, K.N.; Reddy, P.J. Intrusion Detection using Artificial Neural Network. In In 2009, to reduce the original DARPA and KDD problems, Tavallaee et al. ²³ created a new version of KDD called Proceedings of the 2010 Second International Conference on Computing, Communication and NSL-KDD ²³. In this version, the authors removed all redundant records and added new synthetic ones based on Networking Technologies, Karur, India, 29–31 July 2010; pp. 1–7. the correctly labelled records of the original dataset, so that those record types with a lower presence in the original

12/at Stapshad Aa Khid Neichaeste Sichaetznew Aareset-Timevicerveisen Retectione Sestema Bersedvers Loopining v

regeregitem Behavior in a Recent calcanges is intrusion petersized Proceedings of the individual reduction in

sizentamational/W2stksboserRatibas2000e TranlaysadEzaokein2nd tottapett.2000; Debar, H., Mé, L., Wu,

S.F., Eds.; Lecture Notes in Computer Science; Springer: Berlin/Heidelberg, Germany, 2000; pp.

Eversition of the KDD dataset and the application of techniques to rebalance and address consistency

issues, it continues to share the problems of its KDD and DARPA predecessors. Specifically, it relies on 1998 13. Ullah, S.; Ahmad, J.; Khan, M.A.; Alkhammash, E.H.; Hadjouni, M.; Ghadi, Y.Y.; Saeed, F.; network traffic, rendering it outdated in the context of modern network communications and contemporary cyber-Pitropakis, N. A New Intrusion Detection System for the Internet of Things via Deep Convolutional

attacks. Neural Network and Feature Engineering. Sensors 2022, 22, 3607.

12.48 a Ky at a 12,0061+, Dataset. Intrusion Detection in IoT Using Deep Learning. Sensors 2022, 22,

8417.

Given the shortcomings of datasets such as DARPA and KDD with their variants related to the longevity of their 15at Ren 2006, Song Kt, diluss Fruchen da, new dataset called Kyolo 20084, Price vised Intrusion in detection fic from

32 Model Based on Variational Autor Encoder Sensors 2003 to 24 guilt 2009 (almost three years), totalling

100 redtech 83 millive phase reations 035 b. Sin perts duit ap publication vine; and have approximated datas shar and er

a togal en nine ceans of working the solution of 348, including DNS

servers to generate benign traffic. Each record in the dataset provides a total of 24 features associated with the 17. Chandola, V.; Eilertson, E.; Ertoz, L.; Simon, G.; Kumar, V. Minds: Architecture & Design. In Data captured network traffic flows, of which a total of 14 are present in datasets such as DARPA or KDD, while the Warehousing and Data Mining Techniques for Cyber Security; Singhal, A., Ed.; Advances in remaining 10 are new additions, including the labelling of the records, as well as the typology of the detected

attactformatianasecupitypapyingepuBiostenasMAf teSAra2007/ithpne83rea@st historical depth on record, but, in

spite of this, it is still guite balanced. 18. Ahmed, M.; Naser Mahmood, A.; Hu, J. A survey of network anomaly detection techniques. J.

Netw. Comput, Appl. 2016, 60, 19–31. 2.5. Botnet Dataset

19. De Keersmaeker, F.; Cao, Y.; Ndonda, G.K.; Sadre, R. A Survey of Public IoT Datasets for Biglar Beiglar Beiglar

22. Salvatore Stolfo W.F. KDD Cup 1999 Data; UCI Machine Learning Repository, 1999.

23. Tavallaee, M.; Bagheri, E.; Lu, W.; Ghorbani, A.A. A detailed analysis of the KDD CUP 99 data The Cyber Bange Lab at the Autor Metzolog FEEE Symposium on Computational Antelligence for Security and in 2015 using the IXIA Perfect Storm traffic generator 8⁻The simulation environment used to generate the samples consists of three servers, two of which generate benign traffic, while the third is used to generate traffic associated 24/11Figlatus sticks studies D32, explicit state benign traffic. While the third is used to generate traffic associated 24/11Figlatus sticks studies D32, explicit state benign traffic, while the third is used to generate traffic associated 24/11Figlatus sticks studies D32, explicit state benign traffic, while the third is used to generate traffic associated 24/11Figlatus sticks studies D32, explicit state benign traffic, while the third is used to generate traffic associated 24/11Figlatus sticks studies D32, explicit state benign traffic, while the third is used to generate traffic associated 24/11Figlatus sticks studies D32, explicit state benign traffic, while the third is used to generate traffic associated 24/11Figlatus sticks studies D32, explicit state benign traffic, while the third is used to generate traffic associated 24/11Figlatus sticks studies base of the second state of the third is used to generate traffic associated 24/11Figlatus sticks studies base of the second state of the second sta

systems (UNSW-NB15 network data set). In Proceedings of the 2015 Military Communications

2.7. UGR'16 and mormation Systems Conference (MilCIS), Canberra, ACT, Australia, 10–12 November 2015;

pp. 1–6. The UGR'16 dataset ^[26] was created by the University of Granada in 2016 as a result of capturing the real network 26affid actiá-Fredit ám dizz d CBP Catwach 0/a J ch Magán & 2016 n, SBD ; Captor day, Tato dogrob, ePho Filterón JuRy, al G Rúbast, diffetene with also set could be used as a test. The dataset consists of NetFlow traffic flows with almost 17

billion different connections, of which more than 98% were normal traffic, making it very imbalanced. After the traffic 27. Sharafaldin, I.; Habibi Lashkari, A.; Ghorbani, A.A. Toward Generating a New Intrusion Detection was captured, state-of-the-art anomaly detection and network attack identification techniques were employed to tag Dataset and Intrusion Traffic Characterization. In Proceedings of the 4th International Conference the dataset. This involved assigning each record a label indicating the type of attack to which it belonged. Given on Information Systems Security and Privacy, Funchal, Madeira, Portugal, 22–24 January 2018; the size of the dataset and its temporal proximity, it is an updated and current dataset for use in building or training pp. 108–116.
Al and NIDS models.

28. Canadian Institute for Cybersecurity. CSE-CIC-IDS2018. 2018. Available online:

2.8 tto://Datasetsa/cic/datasets/ids-2018.html (accessed on 30 November 2023).

29. Sarhan, M.; Laveghy, S.; Moustafa, N.; Portmann, M. NetFlow Datasets for Machine Learning-The Canadian institute for Cybersecurity (CIC) has generated several datasets to validate the performance of Based Network Intrusion Detection Systems. In Big Data Technologies and Applications. NIDS or to train the models underlying these NIDS. Among the various datasets available, the following should be hig Righted dings of the 10th EAI International Conference, BDTA 2020, and 13th EAI International Conference on Wireless Internet, WiCON 2020, Virtual Event, 11 December 2020; Deze, Z.,

- AlGARS, 241.7 Land in Reperators in Calizmikulati, an symbolic clean work notes of the sum and the comparent of the protection of the sum of
- captured traffic is tagged, and the different attacks that each record corresponds to, including DoS, SSH, and 30. Ring, M.; Wunderlich, S.; Scheuring, D.; Landes, D.; Hotho, A. A Survey of Network-based botnet attacks, are marked in the tag. Intrusion Detection Data Sets. Comput. Secur. 2019, 86, 147–167.
- 31. CAETALS, ICS, 20 Har 1281 a, This Balakausthatic, Antesset connesses to batter a data see the reaction and access requires a prior request (protected). Unlike CICIDS2017, it is
- containing 80 extracted features, and access requires a prior request (protected). Unlike CICIDS2017, it is 32. McHugh, J. Testing Intrusion detection systems: A critique of the 1998 and 1999 DARPA intrusion modifiable and extensible. detection system evaluations as performed by Lincoln Laboratory. ACM Trans. Inf. Syst. Secur.

2.9. NF-UQ-NIDS4

33. Chaabouni, N.; Mosbah, M.; Zemmari, A.; Sauvignac, C.; Faruki, P. Network Intrusion Detection Sarhan et al. [29] have created a synthetic dataset specifically created for machine learning-based NIDS, [29]. This for IoT Security Based on Learning Techniques. IEEE Commun. Surv. Tutor. 2019, 21, 2671– dataset is the result of combining four datasets used in the NIDS domain but transformed into a netflow version. 2701. Two of the datasets used have been analysed previously in this research (UNSW-NB15 [25] and CSE-CIC-IDS2018 341), Sahabithe: MeraphateA-Untrustion Detection and CSE-CIC-IDS2018 341), Sahabithe: MeraphateA-Untrustion Detection and Pathotion and Pathotion and Comparison of the datasets used have been analysed previously in this research (UNSW-NB15 [25] and CSE-CIC-IDS2018 341), Sahabithe: MeraphateA-Untrustion Detection and Pathotion and Pathotion and Comparison of the Automation and CSE-CIC-IDS2018 341), Sahabithe: MeraphateA-Untrustion Detection and Pathotion and Pathotion and Pathotion and CSE-CIC-IDS2018 341), Sahabithe: MeraphateA-Untrustion Detection and Pathotion and Pathotion and Pathotion and Pathotion and Pathotion and Pathotion and Pathotic an

36. Saad, S.; Traore, I.; Ghorbani, A.; Sayed, B.; Zhao, D.; Lu, W.; Felix, J.; Hakimian, P. Detecting

P2P botnets through network behavior analysis and machine learning. In Proceedings of the 2011 **3. Dealing with Labelling Problems in Datasets and the** Minin Annual Hitemational Conference on Phyacy, Security and Hust, Montreal, OC, Canada, 19– **Techniques to Address Them** 21 July 2011, pp. 174–180.

37 La Schlinawi, Apr Schinawi, Whet Travalla e evisori, School and part in School and the sector of the sector of

the different classes that make up the universe being treated. However, when the problem is approached from an 38. Garcia, S.; Grill, M.; Stiborek, J.; Zunino, A. An empirical comparison of botnet detection methods. unsupervised learning perspective, such as anomaly detection, the training dataset is expected to belong to the Comput. Secur. 2014, 45, 100–123. same class. This setup enables the model to learn to identify anomalies by recognizing deviations from the 39 atterns present action of botnet detection from the model to learn to identify anomalies by recognizing deviations from the security present action of botnet detections. The proceeding of the provide the proceeding of the provide the provide

40n Korcoeisstes, offrage on the address of the address of the second second second the state of the address of the as Batmat Datasets in the Interrete of Thingsulige detwork Persion Agady this protein, so at a safe as hive sent tries to detect the noise in the labelling based on loss functions that are insensitive to noise and at the same time 41. Moustafa, N. ToN IOT Datasets; IEEE: Piscataway, NJ, USA, 2019. tries to infer the possible noise in the labelling and in the classification itself [43]. On the other hand, Zhang et al. [44] 42robbsth Chitam Gwark Athaby Rad Aptive Welling No Bercasie Celona bar Net jowshich Tasts Seta Angstahilizen Machierect labeling ning Remeter are in a labeling of the

49:14 Apertmen, J.; Sha, F.; Igel, C. Robust Active Label Correction. In Proceedings of the Twenty-First

International Conference on Artificial Intelligence and Statistics. PMLR, Playa Blanca, Spain, 9–11 When the labelling of the data that make up a dataset is performed manually, there is a risk of unintentional bias April 2018, pp. 308–316. that is intrinsic to the observed data. To address this scenario, a methodology is proposed in ^[45], whose aim is to 4412264796 Jatesheingnating; the Jossible Xidenor gringliferondernut cechladred Quality stering Noisenputational perception tipole FFF Jans deution Netw. Learn. Syst. 2018, 29, 1675-1688.

45. Cabrera, G.F.; Miller, C.J.; Schneider, J. Systematic Labeling Bias: De-biasing Where Everyone is The impact of noise on labeling in artificial intelligence models has also been analysed in several works in a way Wrong. In Proceedings of the 2014 22nd International Conference on Pattern Recognition, that relativises its impact. For example, Natarajan et al. ^[46] propose in ^[46] a simple loss estimator that is unbiased Stockholm, Sweden, 24–28 August 2014; pp. 4417–4422. and minimises the risk of the presence of mislabelled data. Another approach, as proposed by Patrini et al. ^[47], 46 c. Natarajatac Ning Phillogs LeS of Ravikum aboling, Texraciulary Learstende with invision babels us Advances, inclinuNguzeluIntarmetian ReasonsinghSystemsingustagespaciatesclatureBad Harbor Nie USS fantian in instances of 4 hislabelled data [47]. More recent is the work of Wei et al. [48], who this problem and propose two

44 tasats with noise 22 th A langling to Acrive aska the of mark thankers the bow wohes the woodels robterstridues ber to Weisen Are abelling rection Approach. arXiv 2017, arXiv:1609.03683.

- 481 Waricular Zelevance Bengwark bundrin Nill, ea al Linz, Wilier Entalge with Delawirk apelo Energiai test Subserv of 10 Using the cal aborted Human war notations as Xiv. 2012 rockies 11 912 2028 ftv of labelling of the training data.
- 47. in appreach is particularly, internating as the test subsets are assumed to be netlectly taken boards the as the abraic test and the set of test and evaluation meschanism by high the models are tested and validated [42]. Labelling errors in such a dataset
- can destabilise the performance of machine learning models. The datasets tested are those commonly used in the 50. Müller, N.M.: Markert, K. Identifying Mislabeled Instances in Classification Datasets. In field of computational perception (such as MNIST or ImageNet), in the field of language processing (such as IMDB proceedings of the 2019 International Joint Conference on Neural Networks (IJCNN), Budapest, or Amazon Reviews), and finally in the field of audio processing (AudioSet). The results obtained show that there Hungary, 14–19 July 2019; pp. 1–8 are labelling errors that, in some cases, reach up to 10% of the labelling error.
- 51. Hao, D.; Zhang, L.; Sumkin, J.; Mohamed, A.; Wu, S. Inaccurate Labels in Weakly-Supervised Confident Learning (Autionastublidentificacionalearcion betwee an upgreised and comic upervised learning that focuses on characterizing on discrime the Haballing to find 2020, or rect soro resin the labelling in order to train robust
- models. To achieve this, they use data-pruning techniques to clean the dataset before training the models. In [49], a 52. Bekker, A.J.; Goldberger, J. Training deep neural-networks based on unreliable labels. In generalised CL strategy is proposed that is able to find the errors in the labelling by estimating the correct Proceedings of the 2016 IEEE International Conference on Acoustics, Speech and Signal distribution of correct and incorrect labels. Furthermore, it is tested on image datasets, yielding models with higher Processing (ICASSP), Shanghai, China, 20–25 March 2016; pp. 2682–2686.
- 53. Cordero, C.G.; Vasilomanolakis, E.; Wainakh, A.; Mühlhäuser, M.; Nadjm-Tehrani, S. On Generating Network Traffic Datasets with Synthetic Attacks for Intrusion Detection. ACM Trans.

Mül**R**rivan**Betturk20**²⁹, **p2ep3e39t** to detect errors in the labelling of image, text and numerical datasets ^[50]. As a result of the application of this tool, the set of observations of the dataset with a high probability of being 54. Guerra, J.L.; Catania, C.; Veas, E. Datasets are not enough: Challenges in labeling network mislabelled is obtained. This method has been tested on a total of 29 different datasets, both real and synthetic traffic. Comput. Secur. 2022, 120, 102810. and, according to its authors, has been able to find mislabelling in some of them that had not been detected before. 55. Soukup, D.; Tisovčík, P.; Hynek, K.; Čejka, T. Towards Evaluating Quality of Datasets for Network The**Tequiptic biomain**cdm.**Rtateetings: epittie 2021c**

Another methodology in the field of image processing is proposed in ^[52], where the aim is to train a deep learning model with a dataset where there is no confidence in the labelling of the data. To do this, the model adjusts the internal parameters of the neural network while learning the distribution of noise in the labelling and testing it against classical back-propagation models where the goodness of the labelling is assumed.

In the specific area of datasets aimed at addressing cybersecurity or network traffic problems, previous work is more limited, as the generation of these datasets has additional complications with respect to the more general use cases. In ^[53], Cordero et al. ^[53] the problem is reviewed through a comprehensive analysis of various datasets intended for NIDS. The researchers put forth an enhancement to the Intrusion-Detection Dataset Toolkit (ID2T) dataset generation methodology. Subsequently, they evaluate the effectiveness of the proposed ID2T improvement by assessing datasets generated after its application.

The problem of labelling in the field of network traffic is more complex, since it requires specific low-level knowledge of the traffic in order to be able to correctly classify each flow. In ^[54], an analysis of the methods used for labelling this type of dataset, both automatic and manual, is carried out, identifying the weaknesses of each of the techniques along with their advantages and disadvantages.

Finally, to conclude this analysis of the state of the art in dataset quality, in ^[55], an approach to measuring the quality of a network traffic dataset is presented. This quality is used to compare two datasets, to decide if they are equivalent, or if a better quality dataset is found, whether or not it is appropriate to retrain the machine learning models. The proposal for measuring the quality of a dataset is based on the criteria: (i) completeness as the probability that a dataset record can occur in the domain of the machine learning model to be built and (ii) reliability as the probability of occurrence of misclassified or mislabelled data for each possible class. Based on these two criteria, the applicability of a network traffic dataset to a particular problem can be determined.