

AUTOENCODER BASED NETWORK ANOMALY DETECTION

Mukkuh Ganesh^a, Akshay Kumar^b, Pattabiraman V^c

g.mukkuh2017@vitstudent.ac.in

akshaykumar.satheesh2017@vitstudent.ac.in

pattabiraman.v@vit.ac.in

a, b & c Vellore Institute of Technology, Chennai -600127, India

Abstract— Network security is one of the most critical fields of computer science. With the advent of IoT technologies and peer-to-peer networks, the significance of mitigating security threats has never been higher. Network Intrusion Detection Systems are used to monitor the traffic in a network to detect any malicious or anomalous behavior. Anomalous behaviour includes different types of attacks such as Denial of Service (DoS), Probe, Remote-to-Local and User-to-Root. If an attack/anomaly is detected, custom alerts can be sent to the desired personals. In this paper, we explored the effectiveness of various types of Autoencoders in detecting network intrusions. Artificial Neural Networks can parse through vast amounts of data to detect various types of anomalies and classify them accordingly. An autoencoder is a type of artificial neural network which can learn both linear and non-linear representations of the data, and use the learned representations to reconstruct the original data. These hidden representations are different from the ones attained by Principal Component Analysis due to the presence of non-linear activation functions in the network. Reconstruction error (the measure of difference between the original input and the reconstructed input) is generally used to detect anomalies if the autoencoder is trained on normal network data. Here, we compared the performance of 4 different autoencoders on the NLS-KDD dataset to detect attacks in the network. With just reconstruction error, we were able to achieve an accuracy of 89.34% by using a Sparse Deep Denoising Autoencoder.

Keywords— *Network Anomaly Detection; Artificial Neural Network; Autoencoders; Anomaly Detection;*

I. INTRODUCTION

An anomaly or an outlier is the point in the data that are significantly different from the rest of the dataset. In most cases, these anomalies are very few in number and hence have a large class imbalance. Detection of anomalies have a lot of application in a variety of fields, from cyber security to credit card fraud detection.

In the field of network security, it is crucial to identify rogue packets. These packets could lead to network attacks, which could potentially decrease the availability of systems. Millions of dollars can be lost due to website downtimes. Hence, the fast detection of these packets is critical in nature.

Autoencoders are a type of neural network which can learn the encodings of various inputs and learn to reconstruct them from a smaller/larger dimension with minimal error. These Autoencoders can be trained to represent only the normal activities. Hence, the error between the reconstructed inputs and the actual inputs will be high in

the case of anomalous data points. This can be used to detect anomalies in the network.

In this paper, we have taken different types of autoencoders and trained them on the NSL-KDD-CUP network intrusion dataset. The models were trained only on the normal packets and the reconstruction error was used to classify the unseen packets as attack/normal packets. The encoding parts can also be used as a dimensionality reducer and be stacked together with other ML/DL algorithms to get more discriminative models.

II. LITERATURE SURVEY

Network Intrusion Detection system by using the deep learning paradigm known as Self-taught Learning (STL) was proposed by [1]. They created a two-stage sparse autoencoder and trained it on unlabeled data and then used the learned features and passed through a classifier to identify network intrusions. [2] have used the Fuzzy Rough Clustering algorithm to detect network anomalies and gave results comparative to the K-means result making the Fuzzy Rough Clustering algorithm more efficient.

[3] gives a comparative study on the common machine learning algorithms such as random forest, SVM, Naive Bayes etc. that are trained on 41 features. All the features weren't used in this paper. In [4], a variational autoencoder (VAE) is used and it uses the reconstruction probability as the anomaly score to determine the error of the autoencoder. Contrary to most of the autoencoder methods, Variational AutoEncoders are actually generative models. In [5], the authors have combined the variational autoencoder (VAE) with the data that is split task wise while training the model. This model is an effective system to mitigate catastrophic forgetting. In [6], Autoencoder based dimensionality reduction is explored. This can massively reduce the space required to store data, while also building models which can give good latent representations of the same. Principal Component Analysis was explored to detect the anomalies in a network in [7]. This is noteworthy since Autoencoders with linear activation functions have been shown to be similar to PCA [8]. Apart from deep learning, several other methods also exist for network anomaly detection. Ensemble tree algorithms were used for intrusion detection in [9], whereas local adaptive multivariate smoothening was used to reduce the number of false positives by [10].

III. DATASET

The NLS-KDD dataset [11] was used to conduct the analysis for the following experiments. It is the benchmark dataset for intrusion detection based on machine learning approaches. The training data has 1,25,973 data points whereas the test data has 22,543 datapoints. The 38 different attacks were broadly grouped into 4 categories. Namely,

- Denial of Service Attack (DoS)
- User to Root Attack (U2R)
- Remote to Local Attack (R2L)
- Probing Attack

The remaining packets are classified as normal packets. However, the distribution of packets in the training and testing data are not the same to simulate a real-world scenario. It is illustrated by the figures 1 and 2.

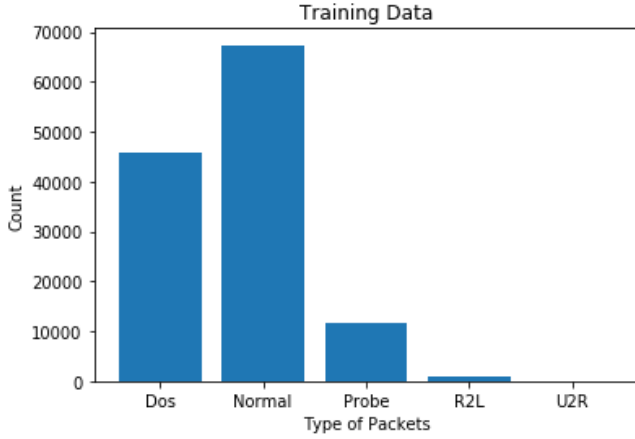


Figure 1: Training Data Distribution.

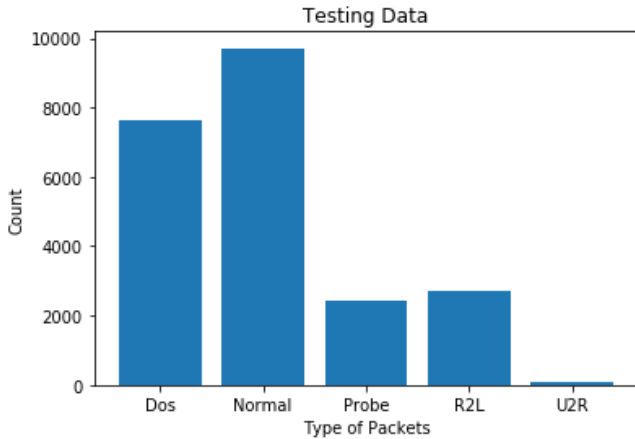


Figure 2: Testing Data Distribution.

The U2R samples are virtually inexistent in the training set (0.04%) and hence, leads to vary poor classification accuracy on trained classifiers. Similarly, the percentage of R2L packets differ significantly in these two datasets. However, by using the Autoencoder method, this hurdle can be overcome, since only the normal packets will be analyzed to build the model.

Label encoding, dummy variables and minmax scaler were used on appropriate columns. The final number of columns in the dataset was 122. 10% of the training dataset was held for validating the models, and the trained models were tested on the test set. For training the autoencoders, only the normal packets from the training data was used.

IV. PREDICTION MODELS

An autoencoder is a neural network that can efficiently learn encodings of data using supervised learning. It consists of two

parts, an encoder and a decoder. The encoder maps the inputs to a new dimension, which could either have higher or lower dimensions than the input dimension, depending on the type of autoencoder architecture. The simplest version comprises of one hidden layer, which maps the input from dimension P to dimension Q. This representation of the inputs is known as the latent representation.

$$g = \sigma(Wx + b) \quad (1)$$

$g \in R^Q$ is the latent representation of the inputs $x \in R^P$, where P and Q are the dimension of the input and latent representations respectively. The second part consists of a decoder which maps the latent representations to the reconstructed inputs.

$$x' = \sigma'(W'g + b') \quad (2)$$

$x' \in R^P$ is the input reconstructed from the latent representation by the decoding layer. The weights of decoding layer could be independent from the encoding layer, or they could also be tied together. The loss for this model is calculated by squaring the difference between the original input and the reconstructed input, and is propagated backwards to the layers.

$$L(x, x') = \|x - x'\|^2 \quad (3)$$

This method of learning is known as self-supervised learning, since no explicit labeled output is needed, and both the inputs and the outputs are essentially the same data points.

A. Undercomplete Deep Autoencoders

An undercomplete autoencoder has lesser neurons than the input layer in its hidden layer(s). When more than one hidden layer is present in the autoencoder architecture, then it is termed as a deep autoencoder. A deep autoencoder can exponentially decrease the amount of data needed to map certain functions, and have experimentally been shown to converge better when compared to shallow autoencoders [12]. An undercomplete autoencoder's structure looks like Fig 3.

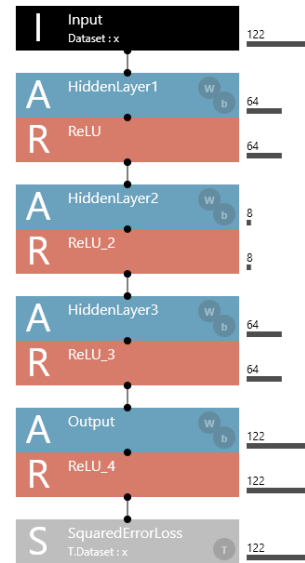


Figure 3: Autoencoder Architecture.

The inputs are reconstructed by the output layer. Loss is calculated by finding the mean of the square of the difference between the reconstructed values and the original values. The loss is then backpropagated through the network and the weights are adjusted accordingly.

By training the autoencoder on just the normal data, one can use the reconstruction error (with set thresholds) to detect the abnormal transactions.

B. Denoising Autoencoder

Its architecture is very similar to that of a vanilla autoencoder. The main difference is that the inputs are corrupted to ensure more generalization, since it forces the neurons/layers to learn more robust features. The type of corruption may vary (e.g., Dropout Layer, adding noise, missing inputs, etc.). During the testing phase, the corruption operation is not performed.

$$\tilde{x} \sim q_D(\tilde{x}|x) \quad (4)$$

x is corrupted into \tilde{x} through a stochastic process. In the proposed architecture, a Dropout layer was added after the Input layer to randomly corrupt/drop the input neurons for every epoch.

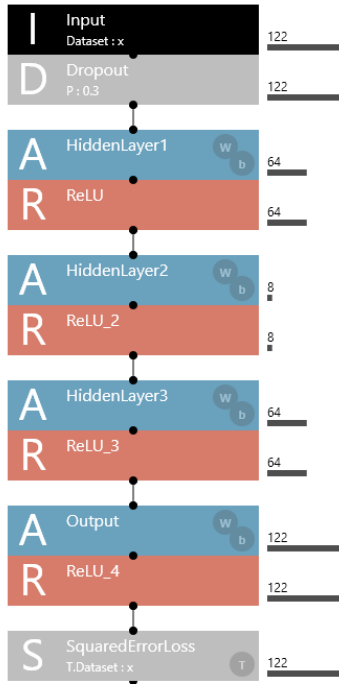


Figure 4: Denoising Autoencoder Architecture.

Here too, the loss is calculated with respect to the input layer and not the corrupted inputs. This has been shown to increase the model's generalization, since it learns to reconstruct the actual inputs from corrupted inputs. Once again, reconstruction error was used to detect abnormal data packets.

C. Sparse Autoencoders

A sparse autoencoder usually contains more neurons in the input layer when compared to the input layer. To prevent the

network from simply copying the inputs into its hidden layer, a sparsity constraint in the form of activity regularizer is added. This is usually L2 regularization.

$$L(x, x') = \|x - x'\|^2 + \text{regularizer} \quad (3)$$

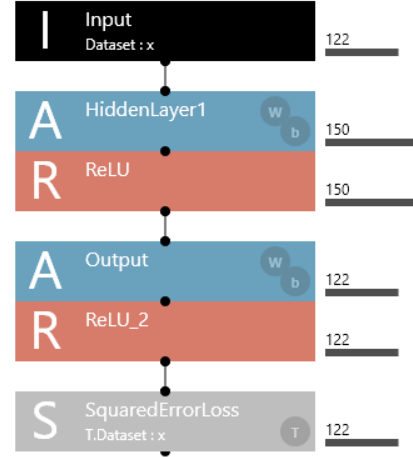


Figure 5: Sparse Autoencoder Architecture.

The model was compiled and trained in a similar manner to the vanilla autoencoder.

D. Sparse Denoising Deep Autoencoder

This is a combination of a Denoising Autoencoder, Sparse Autoencoder and a Deep Autoencoder. A Dropout layer is present after the input layer which corrupts the input, has more neurons in the hidden layer of the encoding part with L2 Activity Regularization.

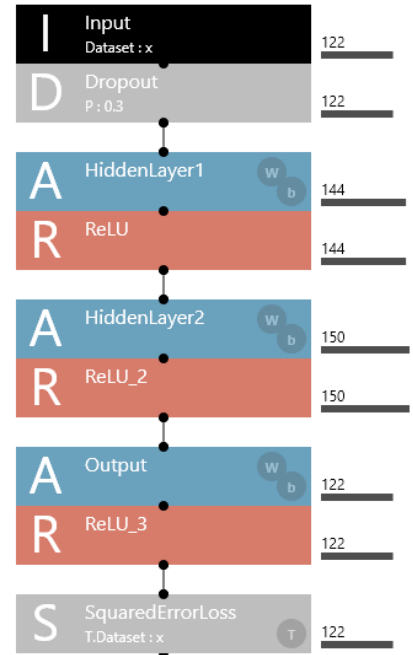


Figure 6: Sparse Denoising Deep Autoencoder Architecture.

The first hidden layer has 144 neurons and the 2nd has 150 neurons. L2 activity regularization of 10e-4 is present in both

these layers. Model is compiled and trained in a similar manner to that of the previous networks.

V. MODEL TRAINING

All the models were trained on the Google Colab platform with the help of its free TPU instances. Tensorflow and Keras were used. The models were compiled with Mean Squared Error as the Loss function, and Nesterov Implemented Adam as the optimizer [13]. Dropout layers in the network had $P = 0.3$. The L2 activity regularizers had a value of $10e-4$. All the models were trained for 25 epochs, during which they converged.

The models converged in 25 epochs. The best model was chosen for evaluation.

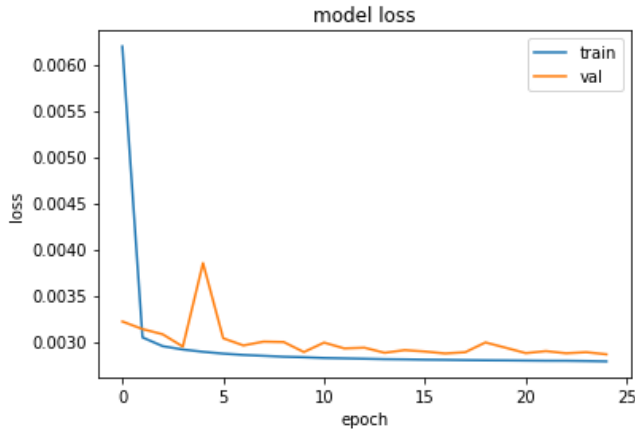


Figure 7: Loss per epoch graph for Deep Autoencoder.

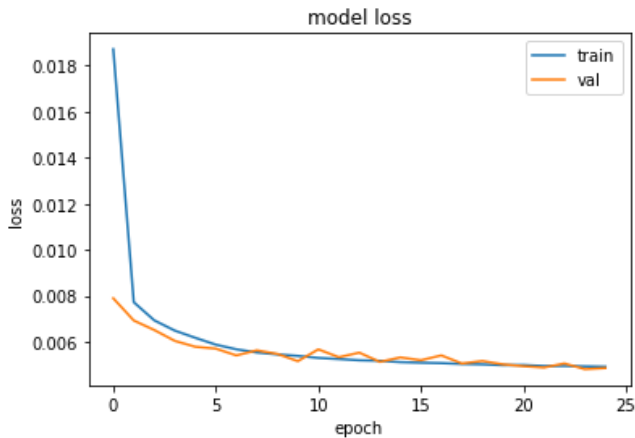


Figure 8: Loss per epoch graph for Denoising Autoencoder.

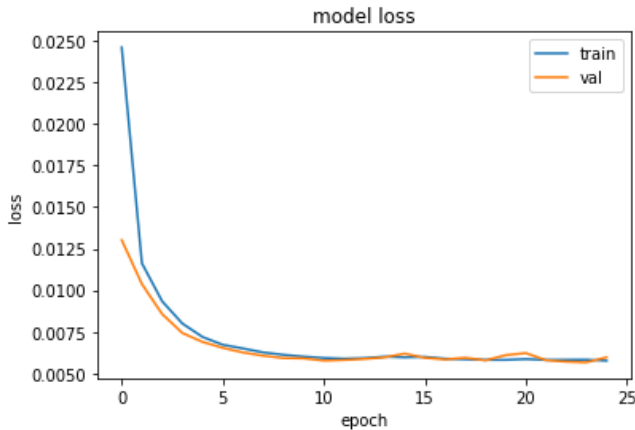


Figure 9: Loss per epoch graph for Sparse Autoencoder.

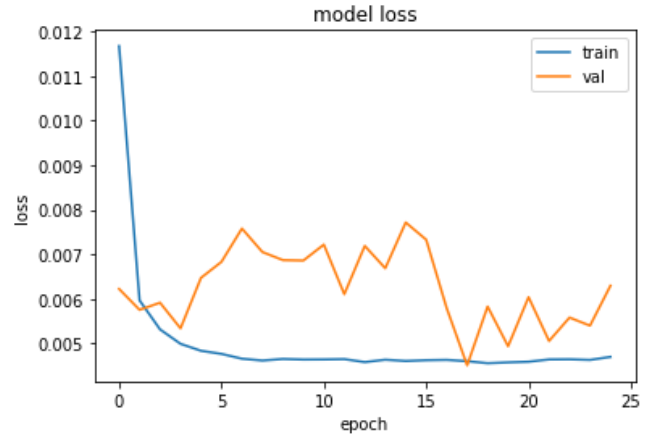


Figure 10: Loss per epoch graph for Deep Sparse Autoencoder.

Table 1: Threshold Values of Autoencoders

Model	Threshold
Deep Autoencoder	0.0028
Denoising Autoencoder	0.0049
Sparse Autoencoder	0.0058
Denoising Sparse Autoencoder	0.0045

The training loss reported by the models in their last epoch is tabulated above. The trained models were then evaluated on the testing dataset.

VI. RESULT DISCUSSION

The thresholds were set as the loss value of the 25th epoch of the respective models. If the reconstruction error of a particular data point is lesser than the threshold, then we classify it as a regular data point. If not, we proceed to classify it as an abnormal datapoint. The results along with their respective violin plots of reconstruction error are plotted below.

Table 2: Confusion Matrix of Deep Autoencoder

Deep Model	Normal Predicted	Attack Predicted
Normal Actual	7219	2491
Attack Actual	715	12118

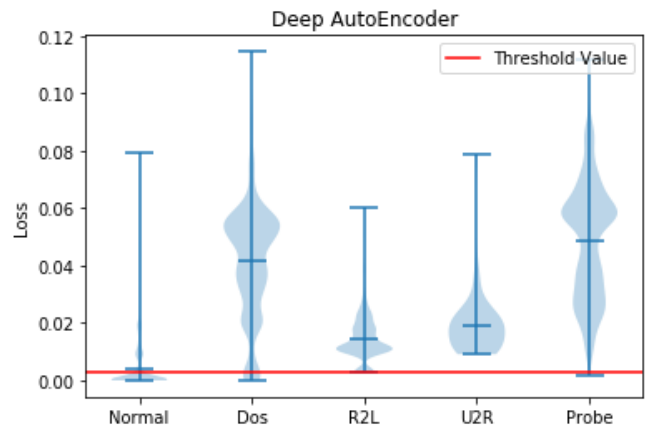


Figure 11: Distribution of Reconstruction error for Deep Autoencoder.

The Deep Autoencoder was able to achieve an accuracy of 85.78% on this dataset. U2R packets were classified as attack

packets with 100% accuracy. The violin plot too shows the U2R reconstruction error to be higher than given threshold value.

Table 3: Confusion Matrix of Denoising Autoencoder

Denoising Model	Normal Predicted	Attack Predicted
Normal Actual	6857	2853
Attack Actual	547	12286

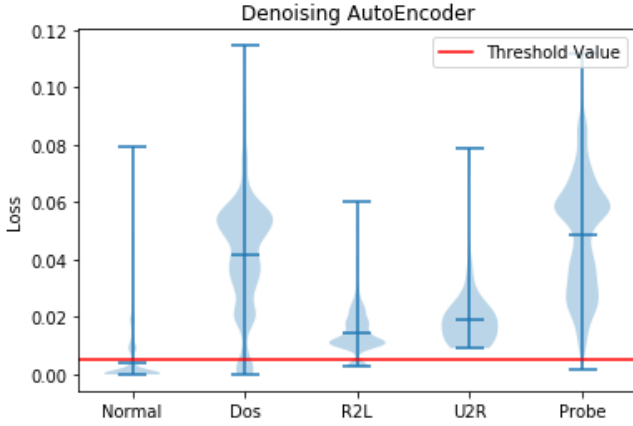


Figure 12: Distribution of Reconstruction error for Denoising Autoencoder.

The Denoising Autoencoder achieved an accuracy of 84.91 % on the full dataset, and reached 100% accuracy in classifying the U2R packets as attack packets. Similar to the deep autoencoder, the violin plot shows that reconstruction error for the U2R packets to be well above the given threshold.

Table 4: Confusion Matrix of Sparse Autoencoder

Sparse Model	Normal Predicted	Attack Predicted
Normal Actual	7025	2685
Attack Actual	636	12197

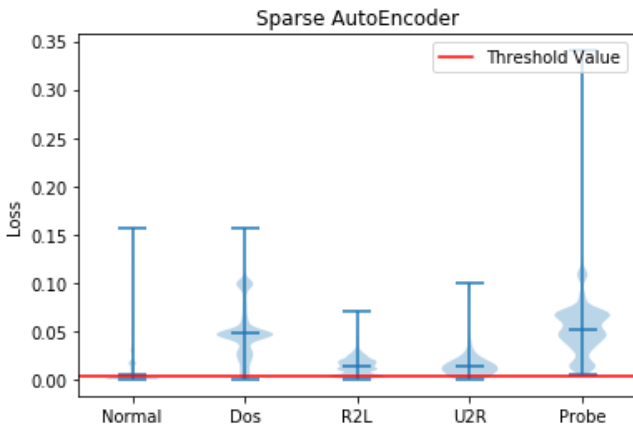


Figure 13: Distribution of Reconstruction error for Sparse Autoencoder.

The Sparse Autoencoder achieves an accuracy of 85.27% on the total dataset and an accuracy of 83.58% on classifying the U2R packets as attack packets. As shown by the violin plot,

some of the packets have their reconstruction error lesser than the set threshold.

Table 5: Confusion Matrix of Sparse Denoising Deep Autoencoder

DSD Model	Normal Predicted	Attack Predicted
Normal Actual	79942	1768
Attack Actual	633	12200

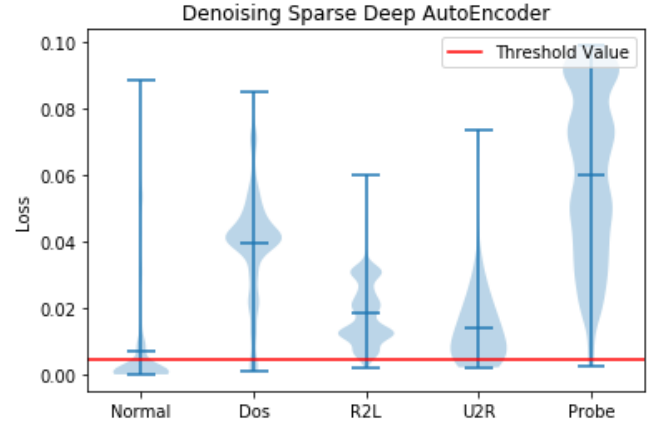


Figure 14: Distribution of Reconstruction error for Denoising Sparse Deep Autoencoder.

The Hybrid Autoencoder achieves an accuracy of 89.34% on the dataset with 94.02% accuracy of U2R classification.

Table 6: Accuracy comparison of Models

Model	Total	Normal	DOS	R2L	U2R	Probe
Deep	85.78	74.35	92.10	87.96	100	99.95
Denoise	84.91	70.62	93.54	98.04	100	99.95
Sparse	85.27	72.35	96.12	87.85	83.58	100
DSD	89.34	81.8	92.58	98.70	94.02	98.84

We can see that the Denoising Sparse Deep Autoencoder outperforms the other types of Autoencoders. It was able to do so with just the reconstruction error. The (virtual) absence of U2R packets in the training data did not hamper the model's performance in classifying them as attack packets. This is one of the main advantages of an Autoencoder based anomaly detection, where the model learns the distribution of a particular type of data and uses it to distinguish the other data types from the same.

VII. CONCLUSION AND FUTURE ENHANCEMENT

The NSL-KDD benchmark dataset was used to train 4 different types of autoencoders. The autoencoders were trained on the normal packets to distinguish them from the attack packets. The combination of 3 models gave us the best results. The Denoising Sparse Deep Autoencoder was able to achieve an accuracy of 89.34% on the unseen test data. Class imbalance due to the presence of extremely low records of U2R attacks in the training set was mitigated by the use of autoencoders. 2 of the models were able to achieve 100% accuracy in classifying U2R type packets as anomalies.

For future research, these models can be saved and the encoding parts can be used as dimensionality reducers. The output of these encoding parts can be used as inputs for other discriminatory ML and DL models such as SVM, Logistic Regression, MLP etc. The use of different sampling techniques such as SMOTE and GMM-based oversampling can also be explored to mitigate the class imbalance for supervised prediction tasks.

VIII. REFERENCES

- [1] Javaid, Ahmad, et al. "A deep learning approach for network intrusion detection system." *Proceedings of the 9th EAI International Conference on Bio-inspired Information and Communications Technologies (formerly BIONETICS). ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering)*, 2016.
- [2] Chimphee, Wittha, et al. "Anomaly-based intrusion detection using fuzzy rough clustering." *2006 International Conference on Hybrid Information Technology*. Vol. 1. IEEE, 2006.
- [3] Revathi, S., and A. Malathi. "A detailed analysis on NSL-KDD dataset using various machine learning techniques for intrusion detection." *International Journal of Engineering Research & Technology (IJERT)* 2.12 (2013): 1848-1853.
- [4] An, Jinwon, and Sungzoon Cho. "Variational autoencoder based anomaly detection using reconstruction probability." *Special Lecture on IE 2.1* (2015).
- [5] Wiewel, Felix, and Bin Yang. "Continual Learning for Anomaly Detection with Variational Autoencoder." *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019.
- [6] Wang, Yasi, Hongxun Yao, and Sicheng Zhao. "Auto-encoder based dimensionality reduction." *Neurocomputing* 184 (2016): 232-242.
- [7] Huang, L., Nguyen, X., Garofalakis, M., Jordan, M. I., Joseph, A., & Taft, N. (2007). In-network PCA and anomaly detection. In *Advances in Neural Information Processing Systems* (pp. 617-624).
- [8] Ladjal, Saïd, Alasdair Newson, and Chi-Hieu Pham. "A PCA-like autoencoder." *arXiv preprint arXiv:1904.01277* (2019).
- [9] Kevric, Jasmin, Samed Jukic, and Abdulhamit Subasi. "An effective combining classifier approach using tree algorithms for network intrusion detection." *Neural Computing and Applications* 28.1 (2017): 1051-1058.
- [10] Grill, M., Pevný, T., & Rehak, M. (2017). Reducing false positives of network anomaly detection by local adaptive multivariate smoothing. *Journal of Computer and System Sciences*, 83(1), 43-57.
- [11] Tavallaee, M., Bagheri, E., Lu, W., & Ghorbani, A. A. (2009, July). A detailed analysis of the KDD CUP 99 data set. In *2009 IEEE Symposium on Computational Intelligence for Security and Defense Applications* (pp. 1-6). IEEE.
- [12] Hinton, G. E.; Salakhutdinov, R.R. (2006-07-28). "Reducing the Dimensionality of Data with Neural Networks". *Science*. 313 (5786): 504–507.
- [13] Dozat, Timothy. "Incorporating nesterov momentum into adam." (2016) (http://cs229.stanford.edu/proj2015/054_report.pdf) last accessed 2020/09/30