

Network Threat Intelligence

December 17, 2023

Shaun Campbell

COSC 6510

1 Introduction

Cybercrime is a growing threat to companies, as business organizations encountered 925 cyberattack attempts per week in 2021¹, and \$10.3 billion was lost due to cybercrime in 2022². Most employees are aware of common attack methods such as phishing and know not to intentionally download untrustworthy items. However, internet users can get viruses just from visiting a website. Exploit kits can reroute users to webpages that scan for vulnerabilities and inject malicious code to install malware on the computer, and exploit kits can even be deployed through ads³. This makes just every-day internet use a source of risk for cyberattack attempts.

Dobberstein Law Firm has a plethora of sensitive information, including personal information on case parties such as social security numbers, names, addresses, and dates of birth, as well as financial information such as credit card numbers, bank account numbers, etc. The nature of the data housed on company servers make it especially at risk to cyberattack attempts.

Although the company has sophisticated cybersecurity infrastructure to guard against attacks, anti-virus software is not 100% infalible⁴. Additionally, firewalls have shortcomings in the sense that they are typically rule-based (e.g., allow connections from a specific source, block connections that use a specific protocol)⁵ and compare data against *known* attack vectors and code⁶. The shortcomings of traditional cybersecurity tools demonstrate the need for another layer of protection, and a way to identify malicious traffic that differs from previously identified threats and avoids firewalls.

The objective of this project is to develop a classification model that can identify malicious traffic using the meta-data of the traffic itself. Malicious traffic can be thought of as suspicious activity on a network. For example, a malicious connection might scan for vulnerabilities in the connecting machine, attempt to inject malicious code, or force an unwanted download.

Internet traffic, at its most granular level, consists of small units of data being sent and received. Data is divided into small units called packets, formatted based on the protocol, then sent across networks. Autonomous System Numbers (ASN) are assigned to owners and operators of IP addresses and can be used to get more information on organizations that own specific IP addresses or ranges. Some definitions of network terms are included in Table 1.

Having a model that can identify malicious traffic will have several benefits for the company. Currently, incident reports can be viewed to see what traffic has been blocked or flagged. However, a model evaluating the traffic will be beneficial for understanding the profile of the internet traffic of the company as whole. It will give insight into what level of threat the company is exposed to due to its internet users and not just known threats that have been avoided. The quantification of malicious exposure will also make possible the evaluation of security policy changes to see how much malicious traffic is avoided or gained due to changes of policy, such as updated block lists and security rules. Finally, it will add a layer of protection in the form of being able to identify large spikes in exposure to malicious traffic.

Term	Definition
byte	Unit of measurement for the size of data sent or received.
packet	Small segment of a larger message. Data is chopped up into packets to be transmitted over networks.
protocol	A set of rules for formatting network data
IP Address	Unique identifier for each computer using the internet
ASN	Autonomous System Number. An ID number assigned to entities that own/operate IP addresses.

Table 1: Definitions¹⁰

2 Business Intelligence Methods

2.1 The Business Process

The project has two key deliverables that will determine its success. Primarily, the end product is a model that can classify network transaction as malicious or benign. As a result, the model will allow additional analysis of the company's threat exposure profile. The IT department will be the users of these deliverables.

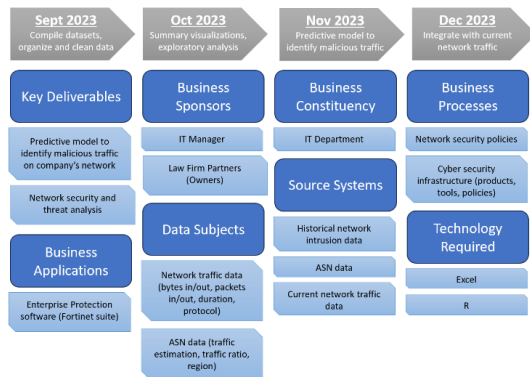


Figure 1: Business Intelligence Roadmap

data. R and Excel are the technologies required for the project. A schematic of the foregoing information as well as an anticipated timeline is displayed in Figure 1.

An examination of the company concludes that it is ready for this BI project. Current data capture capabilities mean the necessary data are available to input into the classification model. The company also has sufficient expertise and experience to not only build the model, but to extract data feeds, input them into the model, and interpret the results. Analytical commitment and openness to organizational change are also present, as the IT department values analysis of internal security and is always striving to improve security. There is no financial commitment needed at this time so that is not of issue.

The anticipated outcomes of the project are three-fold. As outlined, a model to classify network traffic is the first outcome. Specifically, logistic regression will be used because the target variable is binary and it is better suited for classification versus linear probability models, because the latter can result in probabilities less than zero or greater than one, and are more sensitive to outliers which the dataset has many of. Secondly, hypothesis testing will be used to evaluate changes in the predicted malicious traffic exposure levels. Finally, the end outcome is a better understanding of the company's traffic profile, the ability to quantify exposure levels, and the ability to determine the effects of security policy changes.

2.2 Data and Data Structures

The historical labeled network traffic used to build the model is LUFlow Network Intrusion Detection Data Set³, made available via Kaggle. The dataset was compiled using honeypots – fake and unprotected data that would be valuable to hackers, thus attracting attack attempts – at Lancaster University. The network transactions, or flows, are labeled as benign or malicious using cyber threat intelligence sources. Flows undetermined as malicious or benign are classified as “outliers”. Data points include bytes transferred in and out, packets transferred in and out, the duration of the transfer, the protocol used, entropy, and inter-packet arrival time.

The dataset is very large so around 150,000 records are randomly sampled to make up the training data. Due to the abundance of data, omission is used to remove records with missing values. Records that have no bytes sent or received are dropped, since these records do not constitute network flows. Outliers are re-labeled as benign, allowing the model to be built under the premise of classifying truly malicious traffic, rather than classifying as “may-or-may-not-be malicious”. Entropy and inter-packet arrival time are

dropped as variables because these data are not captured currently on the company's network. Indicator variables are created for protocol type 1 and type 17, using type 6 as the baseline.

Plotting histograms of the numeric variables reveals all five are highly left-skewed (Appendix 1). The variables all contain 0 values, so they are transformed using $\log(x+1)$. The resulting distributions tend to still be skewed left but in general are much more workable for further analysis (Appendix 2). Boxplots are generated to check for outliers in the logged variables (Appendix 3). Outliers account for a large percentage of the dataset, so they are retained not to exclude a large subset of data.

Figure 2 displays box plots comparing the logged values of the variables between malicious and benign records. Notable is the difference in bytes out, with benign flows typically having more bytes transfer out than malicious flows. Opposingly, benign transfers appear to typically have less transferred bytes in. Interestingly, several of the variables have less dispersion for malicious traffic signifying that benign traffic may be more variable. This lends evidence that malicious traffic has a different profile than benign, with tighter ranges of values, more bytes transferred in, and less bytes transferred out. This would make sense if a portion of malicious traffic is composed of re-used attack methods attempting to inject malicious content.

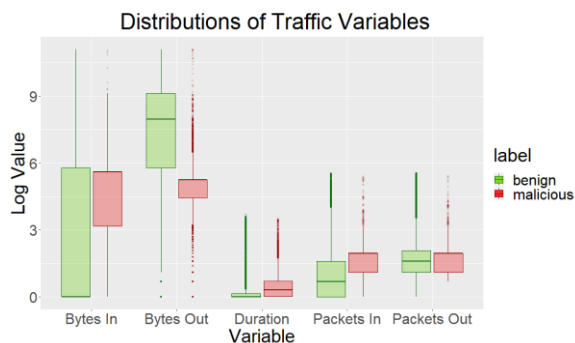


Figure 2: Box Plots for Traffic Variables

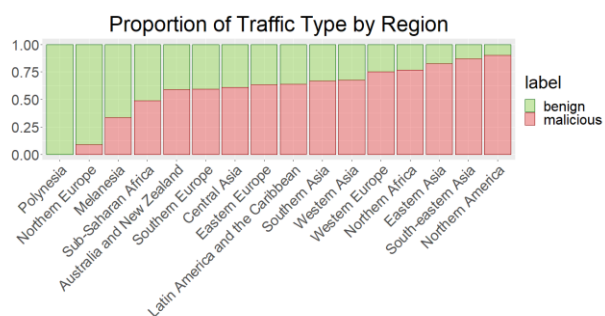


Figure 3: Stacked Bar Chart of Traffic Type by Region

BGPView API⁴ is used to obtain additional data on the ASNs, which includes the country the ASN is based in, the traffic ratio, traffic estimation, and date the ASN was assigned. The R script used to obtain the data is included in Appendix 5. Traffic estimations and traffic ratios for the ASN entities are given as ranges, so indicator variables are created for each range, with the most common range used as the reference category. The country is reduced to region using the UN geoscheme to simplify creation of indicator variables, and North America is used as the reference category. Date allocated is converted to the age, in days, since the ASN was assigned. These data are merged on the source ASN from the network data, therefore providing further data points on the ASN where the network flow originated from.

The resulting dataset contains some missing values in the additional ASN data. Records where the region is unknown are omitted. Missing values in age are filled with the mean, and missing values for traffic estimation and traffic ratio are replaced with the most common level.

Figure 3 shows proportions of malicious traffic depending on the region of origin. It is necessary to mention that the proportions are inflated because the honeypots purposefully attract malicious traffic; it would not be accurate to conclude that 90% of all North American internet traffic is malicious. However, it does demonstrate that certain regions have a higher propensity for initiating attack attempts.

Finally, 25% of the dataset is reserved for testing to avoid any snooping during exploratory analysis.

2.3 Statistical Methodology

Relationships among the variables are investigated to determine whether the variables can be expected to have classification utility. Density plots comparing the numeric variables between malicious and benign subsets are generated to identify ranges where significant differences may occur (Appendix 6),

and the probability of a flow being malicious if the variable's value lies in the range of interest is calculated using Equation 1. The overall probability that a selected flow in the dataset is malicious is 30.76%. The probability analysis shows that there are ranges for each variable that have a probability higher or lower than the baseline value, shown in Table 2. Indicators are considered for the ranges of interest but are not ultimately included because the probabilities are not robust enough and overfitting is a concern.

$$P(\text{malicious} | x \text{ in range}) = \frac{P(\text{malicious} \cap x \text{ in range})}{P(x \text{ in range})}$$

Equation 1: Conditional Probability Formula

Variable	Value Range	$P(\text{malicious})$
log(bytes in)	[5, 6]	55.58%
log(bytes out)	[5, 5.5]	58.61%
log(packets in)	[1.75, 2.25]	57.43%
log(packets out)	[1.75, 2.25]	42.09%
log(duration)	[0, 0.5]	24.58%

Table 2: Malicious Probabilities for Variable Value Ranges

T-tests are performed to investigate differences in mean values for the numeric network variables. Despite the skewness of the variables, the Central Limit Theorem would prove the distribution of sample means are approximately normal so t-tests are appropriate. Results of the t-tests are shown in Table 3. Most of the numeric variables have a significant difference in means across malicious and benign subsets, with the exception of duration. Z-tests are also performed for the categorical variables, and most have significant differences in their subset's proportion of malicious instances versus the overall sample proportion (Appendix 7). This lends evidence that malicious and benign traffics have different profiles and can be classified algorithmically.

Variable	Benign Mean	Malicious Mean	SE	t	df	p
log(bytes in)	2.577	4.269	0.023	-73.242	83250	0.000
log(bytes out)	7.355	4.943	0.016	154.895	83250	0.000
log(packets in)	0.924	1.525	0.008	-77.827	83250	0.000
log(packets out)	1.733	1.598	0.006	21.474	83250	0.000
log(duration)	0.439	0.449	0.007	-1.498	83250	0.134

Table 3: T-Test Results for Network Variables

A pairs plot is generated to determine if there are possible linear relationships amongst the variables (Appendix 8). A relationship of note is packets in and bytes in. Since malicious flows may be attempting to inject code, and bytes of data are compiled into packets, it is feasible that the relationship between the two variables is different from benign flows. A scatter plot, displayed in Figure 4, shows that although the relationship is clearly not linear (changing variance) benign traffic has a prominent subset of points with a strong linear relationship. This relationship is not present for malicious data and the data points have more dispersion. Despite clear nonlinearity, linear regression is performed just to further

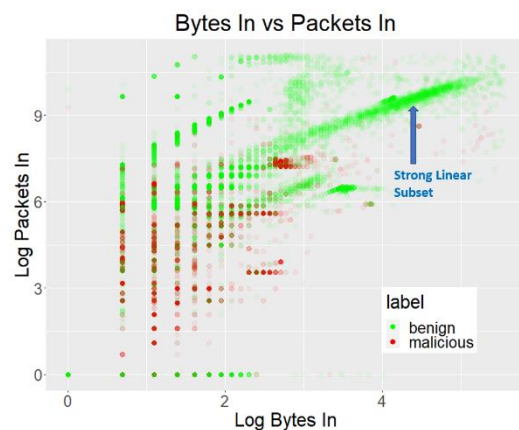


Figure 4: Scatter Plot of Log Bytes In and Log Packets In

investigate the relationship. For malicious traffic the coefficient for log packets in is 2.134 ($R^2 = 0.4796$) and for benign traffic the coefficient is 2.452 ($R^2 = 0.6779$). The higher R^2 indicates more of the variance in bytes in is explained by packets in for benign traffic. The full results of the regression analysis are included in Appendix 9. The assumptions of linear regression are not met so no conclusions can be made, but insight on a possible relationship is gained so an interaction term will be considered for the final model.

Finally, logistic regression is used to build the classification model. Three different models are considered, described in Table 4, and each is evaluated on the training data using 5-fold cross validation. Accuracy is used as the metric to select the best model.

3 Results

The results of the 5-fold cross validation show that the model with the network traffic data, an interaction term between log packets in and log bytes in, and the ASN data performs the best with an accuracy of 86.27% and is selected for the final model. The model with just the base network data has an accuracy of 74.26% in the training set and adding the interaction term increased the accuracy to 78.59%.

Model	Description
Model 1	base network traffic data
Model 2	base network traffic data with packets in * bytes in interaction term
Model 3	base network traffic data with interaction term and ASN data

Table 4: Descriptions of Models Evaluated

The models are used to predict the labels of the training set to obtain performance metrics, which are displayed in Table 5. ROC curves are calculated for the models, and are shown in Figure 5. While the interaction term caused a minimal increase to the area under the ROC curve, the F1 score improved significantly. This means that the interaction term improved the model's precision, recall, and balance between the two. The AUC for Model 3, 0.9324, indicates that the model is quite robust at classifying malicious and benign traffic.

Model	Accuracy	AUC	F1
Model 1 (Train Set)	0.7428	0.8481	0.5628
Model 2 (Train Set)	0.7857	0.8500	0.6775
Model 3 (Train Set)	0.8628	0.9308	0.7742
Model 3 (Test Set)	0.8648	0.9324	0.7773

Table 5: Model Performance Metrics

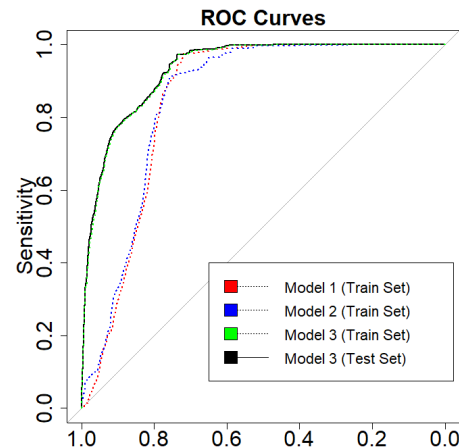


Figure 5: ROC Curves for Evaluated Models

Most of the coefficients in the logistic regression model are significant. The network data variables are all individually significant with the exception of duration. This aligns with the results of the t-test that showed insignificant differences in the sample means of duration. A majority of the ASN indicator variables are also individually significant – 13 of the 14 region indicators, 12 of the 17 traffic estimation indicators, and all of the 5 traffic ratio indicators. Age of the ASN is not significant in predicting malicious network flows. Abridged estimates for Model 3 are displayed in Table 6 and full results are included in Appendix 9.

Variable	Estimate	SE	z	p
(Intercept)	3.024	0.110	27.468	< 2e-16
log(bytes in)	-0.154	0.009	-16.242	< 2e-16
log(packets in)	1.210	0.033	37.228	< 2e-16
log(bytes out)	-0.370	0.007	-55.722	< 2e-16
log(packets out)	-0.305	0.034	-8.867	< 2e-16
log(duration)	-0.004	0.024	-0.164	0.870
proto17	0.881	0.055	16.168	< 2e-16
proto1	1.096	0.238	4.599	0.000
...
log(bytes in):log(packets in)	-0.057	0.005	-11.388	< 2e-16

Table 6: Abridged Estimates for Model 3 Logistic Regression

The model is fit on the test set to evaluate performance on new data. The performance metrics for the test set are also shown in Table 5 and the ROC curve is included in Figure 5. The results show that the accuracy is 86.48% in the test set, slightly higher than the train set, and the ROC curve is almost identical. This is good evidence that overfitting was avoided on the training data. A confusion matrix is generated for the test set and displayed in Figure 6. The confusion matrix shows a balanced number of false positives and false negatives, however since the data is imbalanced this means that the false negative rate is much higher than the false positive rate.

The results of the logistic regression have several implications. First, using just basic network telemetry data and simple logistic regression, malicious traffic can be identified with an accuracy of almost 80% (Model 2). This is confirmation that malicious network traffic has an identifiable profile that differs from normal, benign internet traffic. It is not necessary to inspect the actual content of the data that is being sent and received, but instead use only high-level meta-data on network traffic to identify malicious traffic with a reasonable accuracy.

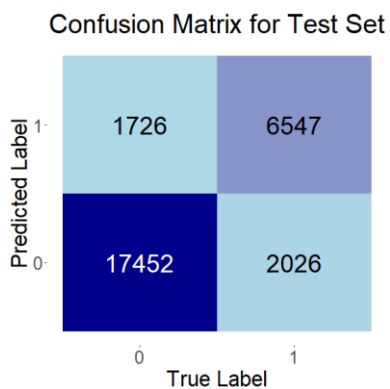


Figure 6: Confusion Matrix for Test Set
1 = malicious, 0 = benign

There is also evidence that the interaction between bytes in and packets in is significant in predicting malicious traffic. The interaction term was individually significant and raised the accuracy and F1 score of the model. The AUC had little improvement, but since the data is imbalanced the F1 score is a better metric. The scatter plot in Figure 4 shows a clear difference in the bytes in/packets in relationship between benign and malicious classes, and the interaction term is likely capturing at least some of this. In other words, there is evidence that the change in likelihood of a transfer being malicious given a change in bytes in is dependent on the value of packets in.

The estimated coefficients give insight on characteristics of malicious traffic. Bytes out and packets out have negative estimates, meaning that an increasing number of bytes and packets transferred out decreases the model's determined likelihood that a flow is malicious. Conversely, packets in is positive, meaning that an increase of packets in increases predicted malicious likelihood. Heavy outbound and mostly outbound traffic ratios are positive, meaning entities who primarily transfer data out appear more likely to initiate malicious transfers. In summary, the more traffic characteristics skew outbound the more likely it is to be benign, and the more they skew inbound the more likely malicious content is being transferred.

Adding information on the organization that owns the foreign IP address associated with the network flow increased performance significantly. This means that obviously, some entities are more likely to initiate malicious traffic. However, characteristics of an entity can be used as indicators without having to know whether the entity itself is a malicious actor. Many security tools use continuously updated lists of compromised or malicious IPs to filter traffic. But qualities of the operating entity such as its region, traffic ratio, and traffic estimation can be used to predict if an entity is prone to malicious activity without having to compare to a predefined, and possibly not comprehensive, list.

More broadly, the final model has utility for the company and in the cybersecurity space. The model uses basic network data that should be trackable with almost all enterprise protection software. Complex metrics such as entropy or inter-packet arrival time are not required. The model also uses publicly accessible information on ASN entities, retrievable with an API. Additionally, the model performs reasonably well without computationally expensive classification algorithms, meaning the model could be a lightweight option for analyzing a company's exposure to malicious traffic. High rates of malicious traffic would signify an organization should consider a change in its policies and protection tools.

To demonstrate use of the model in a real-world scenario, the model is applied to live company traffic from Dobberstein Law Firm. Data for two days, July 28, 2023 (Day 1) and November 27, 2023 (Day 2), are selected arbitrarily to compare. The dates are four months apart, and since security policies and rules are continuously modified, several changes would have taken place in this time period. The live data also only includes successful transfers and thus does not contain any instances blocked by firewalls.

The data is prepared using the same steps as the historical data, with the variables being renamed, merged with ASN data, indicator variables created, and missingness handled. The model is used to classify each record of network traffic as malicious or benign. A total of 144 network transfers were classified as malicious on Day 1 and 49 on Day 2.

Since the number of records for each day is not the same, a Z-test is performed to compare the proportions. The hypothesis test equation is displayed in Equation 2 and the results are shown in Table 7. The sample estimate for Day 1 is 0.178% and the sample estimate for Day 2 is 0.072%. The Z-test rejects the null hypothesis and it can be concluded that the proportion of malicious traffic encountered on the two days is not the same.

$$H_0: \hat{p}_1 = \hat{p}_2 \quad \hat{p}_1: \text{Sample proportion of malicious traffic on day 1}$$

$$H_1: \hat{p}_1 \neq \hat{p}_2 \quad \hat{p}_2: \text{Sample proportion of malicious traffic on day 2}$$

Equation 2: Hypothesis Test for Difference in Proportions

Day	n	\hat{p}	SE	z	df	p	95% CI, Lower	95% CI, Upper
Day 1	81090	0.001776	0.0001801	5.8609	1	2.47E-08	0.0007026	0.001409
Day 2	68035	0.000720						

Table 7: Z-Test Results for Difference of Malicious Traffic Proportions

Further, since the 95% confidence interval does not include zero, it can be reasonably concluded that the proportion of malicious traffic on Day 2 was probably less than Day 1. This is positive information for the company, showing that any policy changes made in the four-month period could have reduced the amount of malicious traffic the company was exposed to. Going forward, the company could make policy changes one at a time and test the difference in malicious traffic exposure to evaluate the effects of the change.

A table showing information on the ASN entities who composed the 193 malicious flows is shown in Table 8. Most of the entities that comprised the malicious instances are hosting companies that provide servers for hosting websites and web services. An interesting entity is Taboola.com Ltd, which accounts for about 44% of the malicious transfers. Taboola is a “native advertising” company that allows advertisers to serve ads on webpages. Some research reveals there have been documented cases of malicious actors deploying “malvertising” through Taboola to redirect users to scam pages⁹. This fact is evidence that the model is working and flagging traffic from an entity that has had cases of serving malicious content in the past. It also demonstrates a use case for the model, since none of the Taboola traffic was blocked by the current firewall. The high number of predicted instances of malicious traffic brings attention to an entity the IT team was previously unaware of, and IT team should consider blocking Taboola IPs.

Description	Malicious Flows	County	Region	Traffic Estimate	Traffic Ratio
Taboola.com Ltd	85	Israel	Western Asia	500-1000Gbps	Not Disclosed
NetActuate, Inc	52	United States	Northern America	10-20Tbps	Mostly Outbound
Akamai International B.V.	18	Netherlands	Western Europe	100+Tbps	Heavy Outbound
OVH SAS	18	France	Western Europe	10-20Tbps	Heavy Outbound
CDN77	10	United Kingdom	Northern Europe	20-50Tbps	Mostly Outbound
Google LLC	5	United States	Northern America	500-1000Gbps	Mostly Outbound
Zayo Bandwidth	2	United States	Northern America	10-20Tbps	Balanced
Amazon.com, Inc.	1	United States	Northern America	500-1000Gbps	Not Disclosed
DigitalOcean, LLC	1	United States	Northern America	1-5Tbps	Mostly Outbound
Limestone Networks, Inc.	1	United States	Northern America	1-5Tbps	Mostly Outbound

Table 8: ASN Entities with Flows Classified as Malicious

4 Conclusions

The statistical analysis of network traffic data and the classification model give insight into characteristics of malicious traffic. Malicious traffic tends to skew inbound, which makes sense given a common form of malicious traffic is injection of malicious content. The significance of the interaction term demonstrates that in regard to bytes in and packets in, there is more going on than just their individual effects, implying bytes of data are compiled differently in malicious traffic. The results also show that a given network transfer's meta-data can be analyzed to determine the probability that it is malicious, although it is far from deterministic. An accuracy of 86% is meaningful, but there are clearly other factors in play besides high-level characteristics.

Still, the classification model has utility in decision making. Avoided and blocked instances are valuable knowledge, but the model gives insight on potentially malicious traffic that got through firewalls. This can be used to evaluate how good current security policies are at stopping malicious traffic. As demonstrated, it can also be used to compare exposure across time periods to evaluate the effect of policy changes.

Another use of the model is for identifying entities that show patterns of suspicious traffic. As outlined, the IT team was not previously aware of the existence of Taboola or its history of malvertising concerns. Taboola is a legitimate ad server, but it seems malicious actors can leverage it to serve malicious content. The model predicted malicious instances originating from Taboola that got through firewalls so the IT team should consider blocking traffic from IP addresses operated by them. No loss in productivity would come from additional ad-blocking so there can only be positive gain in the form of lowered threat exposure from doing so.

There are some weaknesses of the analysis that should be addressed. Entropy and inter-packet arrival time are available in the historical training data but were dropped because the company does not capture these data points. The two metrics could have had predictive power and benefited the model, although their exclusion does keep the model simple. Additionally, the IP addresses in the historical data are masked with the ASN number, which limits the added predictors to high-level information on the owners/operators of the IP addresses rather than information on the actual IP address itself.

The data is also imbalanced, containing about twice as more records of benign traffic than malicious traffic. This means the model becomes biased towards predicting benign labels and results in a high false negative rate, which in the test set is 23.6% versus a false positive rate of 9.0%. This could have been dealt with by resampling the training data to contain an even number of benign and malicious cases or by adjusting the threshold of predicting a malicious class. However not doing the former allows the model to be built under similar conditions to the real world, and doing the latter would require arbitrary cutoffs that could result in overfitting. The F1 score is evaluated due to the imbalanced data and a value of 0.78 indicates a good balance of precision and recall although there is room for improvement.

Another shortcoming is that data was only used for two days to compare malicious traffic exposure levels. The number of malicious instances the company is exposed to can be expected to fluctuate from day to day. The hypothesis test concludes that the sample proportion of malicious traffic on the two days are different but does not offer concrete proof that this is due to policy changes and not just random fluctuations. Still, it demonstrates a use case of the model that can be expanded on by analyzing more days and looking for trends in a time-series manner.

Additionally, the historic data included intra-net activity from Lancaster University as samples for benign transfers. Since Lancaster is located in England, this means a large portion of the benign cases had a region of Northern Europe. This may have inflated the coefficients for the region indicators since North America was used as the reference category. North America had the highest proportion of malicious traffic in the training data and as a result all region indicators were negative. A large portion of live company data is transfers to and from US IPs, which may have skewed predictions toward more positives than in reality. It would be beneficial to look for American-dominant traffic in future analysis.

Lastly, in terms of logistic regression assumptions, one shortcoming is the high prevalence of outliers in the training data. Outliers were retained because removing them would have caused a loss of about 19% of the data. Including the outliers may have skewed the coefficients determined by the logistic regression model and affected model performance. However, in a sense, outliers are what the model is trying to predict, as malicious traffic is different from benign traffic which is presumably fairly homogeneous. Including outliers means the training data mirrors data that the model can expect to see in the real world. Internet traffic is prone to outliers by nature since large portion of internet activity is repetitive, e.g. every Google search would have a very similar profile. Therefore, deviances from the norm are reflected as outliers.

Despite any shortcomings the BI solution presented here can be beneficial for the company and the cyber security space as a whole. There is currently no good way of identifying malicious traffic that got through firewalls because the firewalls themselves are used to identify malicious instances. This solution allows analysis of the leftovers to get a sense of how firewalls are performing and whether allowed traffic has any areas of concern. The profile of allowed traffic can be analyzed to look for trends or spikes and for the evaluation of policy changes.

The model is also simple, easily integrated with existing cyber security infrastructure, and performs reasonably well in terms of accuracy. More advanced machine learning classification models could be used to improve the accuracy and make the utility of the solution even greater. The concept of the solution could open the door to a new area of analysis in the cyber security space, that being additional and ongoing inspection of traffic deemed safe by firewalls.

5 Sources

1. "Check Point Research: Cyber Attacks Increased 50% Year over Year." *Check Point Blog*, Check Point Software, 10 Jan. 2022, blog.checkpoint.com/security/check-point-research-cyber-attacks-increased-50-year-over-year/.
2. United States, Congress, Internet Crime Complaint Center. *Internet Crime Report 2022*, Federal Bureau of Investigation.
3. Ragsdale, Laikin. "Can You Get a Virus from Visiting a Website?" *Sectigo® Official*, Sectigo, www.sectigo.com/resource-library/can-i-get-a-virus-from-opening-a-website. Accessed 9 Dec. 2023.
4. Townsened, Caleb. "Problems with Anti-Virus Software and Alternative Solutions." *United States Cybersecurity Magazine*, United States Cybersecurity Magazine, 24 July 2019, www.uscybersecurity.net/anti-virus-software/.
5. "What Is a Firewall? Definition and Explanation." *Www.Kaspersky.Com*, Kaspersky, 11 July 2023, www.kaspersky.com/resource-center/definitions/firewall.
6. "Security Firewalls: How Do Firewalls Work?" *N-Able*, 9 Apr. 2021, www.n-able.com/blog/how-do-firewalls-work.
7. Mills, Ryan. "LUFlow Network Intrusion Detection Data Set." *Kaggle*, Lancaster University, 6 Dec. 2023, www.kaggle.com/datasets/mryanm/luflow-network-intrusion-detection-data-set.
8. *BGPView API*, bgpview.docs.apiary.io. Accessed 9 Dec. 2023.
9. Segura, Jérôme. "Tech Support Scammers Abuse Native Ad and Content Provider Taboola to Serve Malvertising (Updated)." *Malwarebytes Labs*, Malwarebytes, 27 Sept. 2017, www.malwarebytes.com/blog/news/2017/09/tech-support-scammers-abuse-native-ad-content-provider-taboola-serve-malvertising.
10. "Learning Center Home | Cloudflare." *Learning Center*, Cloudflare, www.cloudflare.com/learning/. Accessed 9 Dec. 2023.