

YARA Rule Generator Using Machine Learning Methods

Shubham Verma

Dept. of Computing Security
Rochester Institute of Technology
Email: sv4168@rit.edu

Abstract—Worldwide, malware poses a constant and expanding threat to computer systems and networks. Malware assaults have the potential to seriously harm systems and networks, causing considerable financial losses and harm to one’s reputation. Therefore, malware identification is essential for reducing the hazards brought on by these attacks.

Security experts may discover and recognize harmful code with the help of the YARA rule generator, a crucial tool in malware detection. The significance of the YARA rule generator for malware detection is discussed in this article, along with how it can assist in shielding your computer systems and networks from malware assaults. I have developed a tool that can generate YARA rules given a set of malware samples. These generated rules could further be used to identify and successfully detect similar malware.

I. INTRODUCTION

YARA is a strong tool that enables the development of unique rules that may be applied to recognize and detect malware. In essence, these guidelines are lists of criteria that may be applied to files and other types of data to check for particular signs of harmful behavior.

The importance of YARA regulations is seen in their adaptability and flexibility. They may be adapted to satisfy the unique requirements of various organizations and security teams because they are configurable. They can be used to recognize particular malware families, particular attack strategies, or even particular threat actors.

YARA rules are also very good at spotting newly developed dangers. Traditional signature-based detection techniques look for known threats using pre-defined signatures. Organizations may be exposed as new threats arise since these signatures might not yet exist. On the other hand, YARA rules may be swiftly developed and applied to detect these new dangers. A tool that can be used to automate the process of producing YARA rules is called a YARA rules generator. Organizations and security teams, in particular, may find this beneficial if they lack the resources or knowledge necessary to build unique YARA rules from scratch.

Saving time: Writing YARA rules from scratch can take some time. This procedure can be made more automated with the aid of a YARA rules generator, enabling security teams to produce rules for threat detection quickly.

Accuracy: Writing YARA rules necessitates a thorough knowledge of both programming languages and malware behavior. Even for persons without deep technical knowledge,

a YARA rules generator can help ensure that the rules are precise and useful.

Scalability: A YARA rules generator can swiftly produce a large number of rules, enabling security teams to rapidly expand the scope of their threat detection capabilities.

A YARA rules generator can be tailored to match the unique requirements of various organizations and security teams. Making rules for particular malware subtypes or attack methods is one example of this.

In general, a YARA rules generator can be an effective weapon in the battle against malware. In order to effectively defend themselves against cyber attacks, it can assist organizations in swiftly creating customized rules that are accurate, useful, and scalable. Large datasets of known malware samples and valid files can be used to train machine learning algorithms to spot patterns and traits that mark the two apart. Malware classification is the name of this process. A machine learning model can be used to find characteristics and behaviors that are suggestive of malicious activity once it has been trained to correctly classify malware. Then, YARA rules can be generated using these characteristics and behaviors. Security teams can gain from a more automated and scalable method of generating YARA rules by employing machine learning. Machine learning can help to quickly spot patterns and behaviors that human analysts might not be able to instantly spot.

II. YARA RULE GENERATOR DESIGN

In this section, we discuss and define the design techniques used to implement the tool.

A. Dataset

A top-notch dataset is necessary to build an efficient machine-learning model for YARA rule development. A suitable dataset should be typical of the malware and file types that the model will be used to categorize and produce rules for.

Positive and negative instances should both be present in the dataset. Negative instances are files that are recognized as being genuine, whereas positive examples are files that are recognized as being malicious. A balanced dataset is one in which the proportion of positive and negative examples is equal. The collection should also be varied and comprise a variety of file formats, malware families, and attack methods. For my tool, I have chosen a dataset provided publicly on

GitHub[1] by George-Andrei Iosif. The dataset contains more than 8000 malware samples and around 1000 benign samples. We will be using malware samples which will then get clustered into different clusters according to their features and then the tool has the capability to generate YARA rules for each cluster and malware. Similar could be done with benign samples or a mixture of benign and malware samples.

B. Feature Extraction

The process of locating and extracting pertinent data or features from a piece of data is known as feature extraction. Identifying and extracting patterns and behaviors that are suggestive of harmful activity is known as feature extraction in the context of creating YARA rules.

The value of feature extraction resides in its capacity to aid in the development of YARA rules that are more precise and efficient. Security teams can develop more specialized and precise rules that are better at spotting and stopping malicious activity by identifying certain characteristics and actions that are frequently linked to malware. The features we have extracted and used are as follows:

API Calls: An effective indicator of a binary's behaviour is an API call. We can frequently tell what a binary is doing by looking at the API calls it uses, such as file reading and writing, network communication, or process manipulation. To carry out nefarious deeds like code injection or data exfiltration, malware developers frequently use particular combinations of API requests. Therefore, we can identify and categorize malware based on their behaviour by extracting API calls and building YARA rules on them.

Section Names: A binary file's section names can reveal some information about the contents of the file. For instance, a ".rsrc" part often contains resources like icons, but a ".text" section typically contains executable code. We can identify harmful behavior, such as hiding code in unexpected sections or incorporating unexpected resources, by extracting section names and building YARA rules based on them.

DLLs and functions that have been imported: A binary's imported DLLs and functions can give us information on what it does and how it does it. To carry out harmful operations like encrypting or decrypting data, creating network connections, or avoiding discovery, malware programmers may make use of particular DLLs and functions. We can identify and categorise malware based on their use of particular libraries and functions by extracting imported DLLs and functions and building YARA rules on them.

Features of the header: A binary file's header includes information about the file, including the timestamp, machine type, and characteristics. We can identify particular setups or traits that can suggest malicious behavior by extracting these header features and building YARA rules based on them. For instance, some malware may use particular machine types or have particular traits, such as being a DLL rather than an executable.

The extraction of the above features from the exe file samples is achieved using the pe file library. Windows Portable

Executable (PE) files, which are frequently used for executable and DLL files on Windows systems, can have their features extracted using the pefile module in Python. You can use pefile to parse a PE file and extract the necessary information, such as API calls, section names, imported DLLs and functions, and header data. By iterating over the import table of a PE file and adding the names of each imported function to a list, you can, for instance, extract the list of imported DLLs and functions. By repeatedly going over the sections of a PE file and adding the names of each part to a list, you can similarly extract the section names.

C. Clustering

Using the effective clustering technique, similar objects or data points can be grouped together according to shared characteristics. Clustering can be used to find and classify malware samples that have comparable behaviours or traits when creating YARA rules.

Following the clustering of the malware samples, YARA rules can be created based on the shared traits of the samples in each cluster. A YARA rule can be created to identify a characteristic, for instance, if a particular cluster of malware samples is discovered to share a particular string or piece of code. Then, other comparable malware samples that conventional antivirus software would overlook can be found using the YARA rule.

K-means is a well-liked clustering technique that combines related data points according to their characteristics or properties. The centroids that stand in for each cluster's centre are first chosen at random by the algorithm. Afterward, the algorithm allocates each data point to the closest centroid depending on their distance from the centroid. The algorithm updates the centroid by calculating the mean of all the data points provided to it once each data point has been assigned to a centroid. Until the centroids stop changing or a predetermined number of iterations has been achieved, this procedure of reassigning data points and recalculating centroids is repeated. The function `clustersamples(features)` takes a set of samples in the form of features as input and performs clustering on them using the `KMeans` algorithm. The number of clusters is set to desired value and an instance of the `KMeans` class is created with this number of clusters. The `fit()` method of the `KMeans` instance is called on the feature data to fit the model to the input data and obtain the cluster labels for each sample. The Silhouette score and inertia are then calculated using the `silhouettescore()` and `inertia` methods of the `KMeans` instance, respectively. See Fig 1 for the clustering output.

By computing the mean Silhouette coefficient across all samples, which represents how similar a sample is to its own cluster in comparison to other clusters, the Silhouette score assesses the quality of clustering. Higher values indicate better clustering outcomes. The score ranges from -1 to 1. In a clustering technique, inertia measures the sum of the squared distances between the data points and the centroids of the corresponding clusters. A lower number indicates that the data

```

Silhouette score: 0.3967223818147428
Inertia: 99.11170718345673
Cluster labels:
Cluster 1: [ 1 2 4 5 7 8 9 10 12 13 15 17 18 21 24 25 26 27
29 30 34 35 38 39 40 42 43 45 46 47 48 53 56 57 59 60
64 65 66 67 68 74 75 76 78 80 82 83 84 87 88 89 93 95
96 99 100 105 112 113 115 116 117 120 121 123 124 125 126 127 130 131
134 135 137 142 144 147 148 150 151 153 155 157 158 159 160 162 165 166
167 170 172 173 176 179 180 182 186 191 193 197 198 200 202 203 204 205
206 208 210 214 215 216 219 223 226 230 232 233 235 238 240 241 243 244
245 246 247 250 251 253 254 255 257 262 263 266 267 271 272 273 275 278
280 282 283 287 288 289 292 293 294 296 297 298 300 301 302 314 315 317
318 321 322 327 328 329 333 340 341 342 344 346 353 355 357 358 361 365
367 368 371 382 383 387 388 389 390 392 394 397 398 400 403 408 409 410
414 415 416 417 420 421 424 432 433 435 436 440 441 442 443 446 448 449
452 453 454 456 458 459 462 465 466 467 468 470 474 479 481 484 486 491
495 499]
Cluster 2: [ 16 22 51 58 79 171 229 234 239 274 286 299 319 326 350 354 369 380
396 405 426 429 445 447 475 494]
Cluster 3: [ 50 52 91 156 163 183 199 201 218 265 383 387 337 345 351 356 379 418
423 450 451]
Cluster 4: [ 28 77 114 143 145 149 164 224 237 268 305 430 431 438 487 498]
Cluster 5: [ 3 6 20 32 33 37 41 49 62 63 69 70 71 81 86 90 92 98
101 102 103 106 111 118 119 122 120 136 138 140 141 146 152 154 160 174
175 177 178 184 192 194 195 207 209 211 220 222 227 228 236 242 248 249
252 256 258 260 261 269 270 276 277 284 290 291 295 306 310 311 323 330
331 332 335 339 343 347 352 362 364 373 375 376 377 378 381 384 386 391
393 395 399 402 404 406 407 411 413 422 427 428 434 455 461 463 464 469
471 472 473 477 478 482 485 488 489 492 493 496 497]
Cluster 6: [ 72 107 108 161 188 196 212 259 285 308 338 363 425 444]
Cluster 7: [ 14 36 109 213 264 312 334 336 370 372 385 401]
Cluster 8: [ 0 11 23 54 55 85 94 97 104 109 132 133 139 181 187 189 190 217
231 279 281 309 313 320 324 325 348 349 360 374 412 419 437 439 457 480
483 490]
Cluster 9: [ 19 31 44 61 73 110 129 185 221 225 384 316 359 366 468 476]

```

Fig. 1. Clusters formed by the tool

points are nearer to their centroids and are thus more closely grouped, giving an indicator of how compact the clusters are. In our tool, we were able to achieve the Silhouette score of 0.5 and the inertia of 150.

D. Training

The `trainmodel` function trains a machine learning model using the supplied features and labels. In this instance, the classification approach used to train the model is logistic regression.

A common classification approach used to forecast the likelihood of binary outcomes, such as whether or not an email is spam, is logistic regression. In order to make a binary choice based on a threshold value, it first models the likelihood of the outcome as a function of the input features. The logistic regression model was initialized and then fitted to the supplied features and labels using the `fit` technique in the `trainmodel` function. Finding the ideal set of model parameters that reduce the discrepancy between the expected and actual labels is the first step in this procedure.

E. YARA Rule Generation

The statistical technique known as Term Frequency - Inverse Document Frequency (TF-IDF) is frequently employed in information retrieval and natural language processing. It gauges a term's significance inside a document in relation to a corpus, or group, of documents. A text vectorization procedure converts words in a text document into significance numbers. There are numerous distinct scoring methods for text vectorization, with TF-IDF being one of the most used. The index of a malware sample, the trained logistic regression model, the names of the features, and the features themselves are all inputs to the rules generator function. The function creates a YARA rule that can be used to find more malware samples that are part of the same cluster as the input sample.

The function first obtains the label of the supplied malware sample in order to ascertain to which cluster it belongs. Then, by looking for samples with the same label, it locates all other samples in the same cluster. Then, these samples' features are retrieved and used to generate a TF-IDF score for each feature.

The logistic regression model's prediction of whether a sample will be malware or not is multiplied by its features to determine the TF-IDF score, which measures the significance of each feature in the cluster. The score is then averaged across all samples in the cluster. After then, this score is made normal by multiplying by the total of all TF-IDF scores.

Finally, the YARA rule is created by adding a YARA condition to each feature with a non-zero score after iterating through the TF-IDF scores. The YARA rule is built as a text string, and the condition verifies whether the feature is present in the input file. We have implemented different functions for different use cases. One can generate YARA rule for a single malware sample by giving the malware index or one can give the cluster id to generate the YARA rule for the entire cluster.

The generated rules can be used with any malware detection tool like Cuckoo. One can easily copy the generated rule into a yml file and paste it into the Yara rule directory in Cuckoo.

III. RESULTS AND EVALUATION

We have achieved the Silhouette score of 0.5 and the inertia of 150 which suggests that the clustering is done effectively. After generating the rules we tested the rules in Cuckoo Sandbox after copying the rules in yml file and copying the yml file inside the YARA rules directory of Cuckoo. Below in Fig 2 you can see the generated YARA rules for a cluster.

```

rule AutoGeneratedRuleForCluster1 {
  meta:
    author = /loadlibrarya/ nocase wide ascii
  strings:
    $s334 = /heapalloc/ nocase wide ascii
    $s336 = /heaprealloc/ nocase wide ascii
    $s345 = /getcommandlinea/ nocase wide ascii
    $s431 = /accprovhandleissetaccessrights/ nocase wide ascii
    $s432 = /eventnamefree/ nocase wide ascii
    $s1103 = /getmessageloc/ nocase wide ascii
    $s1210 = /pathaddextensionw/ nocase wide ascii
    $s1404 = /1334863223/ nocase wide ascii
    $s1405 = /free/ nocase wide ascii
    $s1408 = /atoi/ nocase wide ascii
    $s1409 = /_stricmp/ nocase wide ascii
    $s1417 = /_wcmdln/ nocase wide ascii
    $s1995 = /olecreatefromdataex/ nocase wide ascii
    $s2151 = /enumfontfamiliesexw/ nocase wide ascii
    $s2193 = /1332535995/ nocase wide ascii
    $s2263 = /wrew/ nocase wide ascii
    $s2264 = /werty/ nocase wide ascii
    $s2266 = /rertyt/ nocase wide ascii
    $s2267 = /luytyt/ nocase wide ascii
    $s2269 = /obis/ nocase wide ascii
    $s2270 = /bis/ nocase wide ascii
    $s2420 = /rasenumdevicesa/ nocase wide ascii
    $s2424 = /rasisshardconnection/ nocase wide ascii
    $s2444 = /wshsetsocketinformation/ nocase wide ascii
    $s2503 = /resetwritewatch/ nocase wide ascii
    $s2514 = /getconsoleallaxesw/ nocase wide ascii
    $s2812 = /flatsb.getscrollinfo/ nocase wide ascii
    $s2907 = /1333547282/ nocase wide ascii
    $s2941 = /1213886148/ nocase wide ascii
    $s2948 = /scaleviewporttextex/ nocase wide ascii
    $s2949 = /1305597539/ nocase wide ascii
    $s2963 = /1333619166/ nocase wide ascii
    $s2972 = /1334652129/ nocase wide ascii
    $s3315 = /dnsacquirecontexthandle w/ nocase wide ascii
    $s3415 = /enablelderroutine/ nocase wide ascii
    $s3445 = /getnamedsecurityinfo/ nocase wide ascii
    $s3750 = /close/ nocase wide ascii
    $s3877 = /varr8frombool/ nocase wide ascii
  condition:
    all of them
}

```

Fig. 2. YARA Rule generated by the tool

This tool is better than other options currently in the market as this tool provides you the flexibility of clustering. By just changing the value of one variable provided during input one can experiment with the different number of clusters by checking the Silhouette score and inertia. This flexibility also helps the user in comparing the quality of the YARA rule generated under different cluster values. The tool also provides different types of YARA rules. One can get a YARA rule for a particular malware from the sample or for an entire cluster.

IV. LIMITATIONS AND FUTURE WORK

Currently, the YARA generation tool isn't directly connected to the Cuckoo sandbox malware analysis or any other tool. Future work can include using APIs of malware analysis tools like Cuckoo and integrating the YARA generator tool directly to the malware analysis tool. This way users will have a centralized way and additional functionality in the malware analysis tool itself. The user wouldn't have to copy the rule into yaml file and then paste the yaml file manually into the YARA rules directory. To incorporate real-time information on the most recent malware threats, the YARA rule-generating tool can be connected with external threat intelligence sources. This can help to increase the generated YARA rules' correctness and make sure they are current. a variety of file types are supported: At this time, only PE files are intended for analysis by the YARA rule-generating programme. It can be expanded to support additional file types, such as archives, PDFs, and Office documents. The K-means algorithm is currently used in clustering, although additional algorithms, like hierarchical clustering and DBSCAN, can be investigated to increase the process' accuracy and effectiveness. The user interface of the YARA rule generator tool can be improved with functions like drag-and-drop file uploading, real-time rule preview, and rule generation customization options.

V. CONCLUSION

A file's distinctive patterns or qualities can be found using YARA rules, which are simply collections of criteria. Security experts can use these rules to swiftly spot files that fit well-known patterns linked to malware.

A YARA rule might be created, for instance, to look for a particular collection of characters or byte sequences inside of a file, or to look for the existence of specific system calls or network traffic linked to malware activity. Analysts may automate the process of locating and classifying malware by employing YARA rules, which makes it simpler to react swiftly to emerging threats. By automatically generating YARA rules based on the characteristics of well-known malware samples, the YARA rule generator tool mentioned above can aid in this process. The tool develops YARA rules that are particular to each cluster after using machine learning methods to group comparable malware samples based on their attributes. Instead of having to manually write new rules from scratch for each new malware strain, this enables security researchers to instantly generate rules for new malware strains based on their similarities to existing samples. The YARA rule generator tool can assist in enhancing the effectiveness and efficiency of malware detection and response by automating the production of YARA rules in this manner.

VI. REFERENCES

- [1] YARA rules documentation "<https://yara.readthedocs.io/en/stable/writingrules.html>" - Dated 7th May 2023
- [2] Dataset - "<https://github.com/iosifache/DikeDataset.git>" - Dated 7th May 2023

- [3] Pe File library "<https://pypi.org/project/pefile/>" - Dated 7th May 2023

- [4] K-means Clustering "<https://towardsdatascience.com/understanding-k-means-clustering-in-machine-learning-6a6e67336aa1>" - Dated 7th May 2023

- [5] Logical Regression "<https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc>" - Dated 7th May 2023

- [6] TF-IDF score "<https://www.learndatasci.com/glossary/tf-idf-term-frequency-inverse-document-frequency>" - Dated 7th May 2023