



SEI-SICITE 2021

Pesquisa e Extensão para um mundo em transformação

Uso de Inteligência Computacional para detecção de eventos anômalos de redes

Use of Computational Intelligence to detect anomalous network events

Guilherme de Almeida do Carmo*, Luiz Fernando Carvalho†

RESUMO

Com a massificação dos recursos computacionais e o surgimento de serviços cada vez mais sofisticados, o gerenciamento da rede vem se tornando uma atividade cada vez mais complexa, por conta de eventos anômalos nas redes. Essa pesquisa tem como o objetivo detectar ataques de Negação de Serviço Distribuída em uma base de dados disponibilizada pela University of New Brunswick, com uso de Aprendizado de Máquina Não-Supervisionado. Utilizando algoritmos dos tipos Hierárquico, Baseado em Erro Quadrático e Baseado em Densidade, concluiu-se que para o problema em questão, os que obtiveram os melhores resultados foram os dois primeiros, enquanto para o terceiro, houve dificuldade na identificação dos ataques.

Palavras-chave: Negação de Serviço Distribuída, Aprendizado de Máquina Não Supervisionado, Hierárquico, Erro Quadrático, Densidade.

ABSTRACT

With the massification of computational resources and the emergence of increasingly sophisticated services, network management has become an increasingly complex activity, because of anomalous events in networks. This research aims to detect Distributed Denial of Service attacks in a database provided by the University of New Brunswick, using Unsupervised Machine Learning. Using three kinds of algorithms: Hierarchical, Based on the Quadratic Error Metric and Density-Based, it was concluded that for the problem in question, the ones that obtained the best results were the first two, while for the third one, there was some difficulty in identifying the attacks.

Keywords: Distributed Denial of Service, Unsupervised Machine Learning, Hierarchical, Quadratic Error Metric, Density-Based.

1 INTRODUÇÃO

A área de segurança da informação tem crescido muito, por conta de programas maliciosos e ataques cada vez mais frequentes nas redes de comunicação. Os ataques que o presente projeto explora são de Negação de Serviço Distribuída (DDoS, do inglês Distributed Denial of Service), sendo uma variação dos ataques de Negação de Serviço. Esse ataque utiliza uma “rede zumbi” de computadores contaminados, que sobrecarrega o servidor da vítima, causando lentidão ou queda total do mesmo.

* Engenharia de Computação, Universidade Tecnológica Federal do Paraná, Apucarana, Paraná, Brasil; guilhermecarmo@alunos.utfpr.edu.br

† Universidade Tecnológica Federal do Paraná, Campus Apucarana; luizfcarvalho@utfpr.edu.br



Existem no ramo de aprendizado de máquina dois tipos de problemas: Classificação e Regressão. Problemas de classificação têm como saída um valor discreto, o qual representa uma das categorias presentes na base de dados. Já os problemas de regressão consistem numa saída de valor contínuo, por exemplo, a altura prevista de um grupo de pessoas. No caso da detecção de ataques DDoS em uma rede, trata-se de um problema de classificação, sendo a saída um valor que indique a presença de um ataque ou não.

Considerando o aprendizado de máquina, existem duas vertentes mais conhecidas: Supervisionado e Não-Supervisionado. A principal diferença entre esses métodos está na natureza dos dados de entrada, já que no supervisionado existem instâncias. Neste projeto, é utilizado aprendizado de máquina não-supervisionado para detecção das anomalias na rede por meio da descoberta de padrões, uma vez que pergunta-se: É possível encontrar um padrão de comportamento dos ataques DDoS no tráfego da rede?

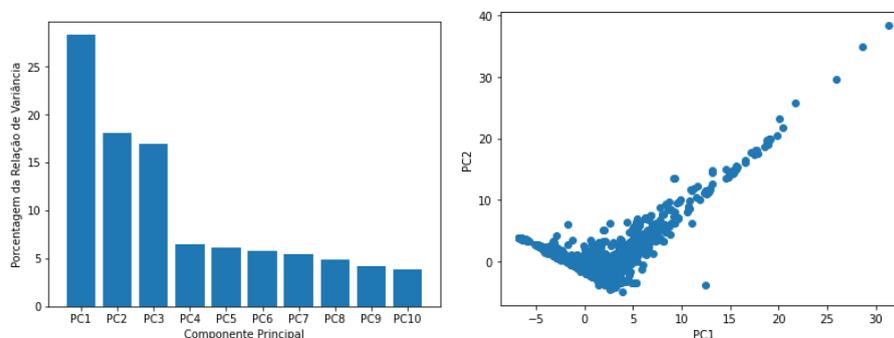
Para isso, deve-se analisar na base de dados cada momento que ocorre um ataque na rede e encontrar os padrões que esse ataque DDoS segue, podendo assim, prever novas ocorrências. Uma forma de analisar os dados é com a detecção de clusters e outliers. “Cluster”, em inglês, significa “grupo”, portanto, clusterização é uma técnica que busca padronizar os dados não rotulados em grupos, quando há muitas dimensões (A. Patcha, J.M. Park, 2007). Já “outliers” são os ruídos (ou anomalias) presentes entre os dados.

2 MÉTODO (OU PROCEDIMENTOS OPERACIONAIS DA PESQUISA)

A base de dados CIC-DDoS, fornecida pela University of New Brunswick (UNB), contém colunas que indicam informações de redes de computadores, como IP de origem e destino, portas de origem e destino, entre outras. Nela, foram executados 12 diferentes tipos de ataques DDoS no dia de treinamento e 7 tipos no dia de teste.

Fazendo uma análise inicial da base de dados, nota-se que a visualização do problema se torna muito difícil, pelo fato de possuir muitas dimensões. Em casos como o em contexto, pode-se utilizar a Análise dos Componentes Principais (ACP). A ACP é uma técnica matemática que permite a redução da dimensionalidade dos dados a partir da determinação dos autovalores da matriz de covariâncias e dos autovetores a eles associados, tornando possível a visualização dos dados em um gráfico com duas ou três dimensões, por exemplo.

Figura 1 – Autovalores considerando os 10 Componentes Principais (à esquerda) e Gráfico da ACP (à direita)



Fonte: autoria própria (2021)

A Figura 1 é uma demonstração de que, a partir dos componentes principais, uma plotagem de gráfico para visualização dos dados passa a ser possível, com os componentes PC1, PC2 e PC3 sendo os eixos x, y e z



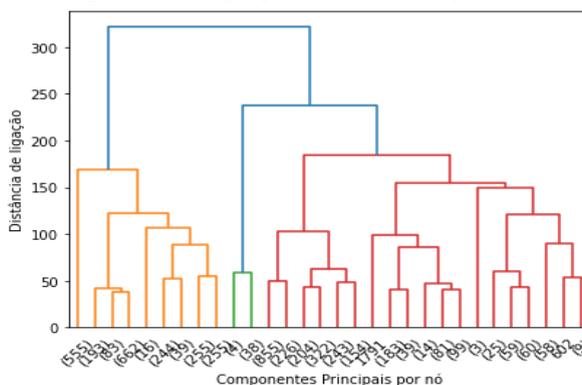
z, respectivamente. Com a geração do gráfico da ACP, o problema torna-se mais inteligível e o modelo passa pela etapa de treinamento e previsão.

Com os dados visualizáveis, o problema passa a ser a detecção dos padrões que cada dado segue na rede. Para isso, foram utilizados três diferentes algoritmos de aprendizado de máquina não-supervisionado: Agrupamento Hierárquico com Abordagem Aglomerativa, K-Means (ou K-Médias) e DBSCAN.

2.1 Agrupamento Hierárquico com Abordagem Aglomerativa

Para a definição do número de clusters a serem utilizados num problema de Agrupamento Hierárquico, pode-se utilizar um dendrograma como auxílio. Um dendrograma é um diagrama dos agrupamentos. Graficamente, ele pode ser representado como na Figura 2, onde nota-se a existência de 3 clusters representados em amarelo, verde e vermelho.

Figura 2 – Dendrograma do Agrupamento Hierárquico com método Ward



Fonte: autoria própria (2021)

Após a definição do número de clusters para o problema, é possível a visualização do gráfico da ACP (Figura 1 à direita) com os dados de características semelhantes agrupados. Para isso, deve-se utilizar uma métrica de integração na execução do algoritmo. Neste projeto, utilizou-se a distância Euclidiana, Eq. (1), sendo x_i e x_j duas instâncias definidas e n sendo os n atributos que descrevem as instâncias x_i e x_j , respectivamente.

$$distancia = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (1)$$

2.2 K-Means (K-Médias)

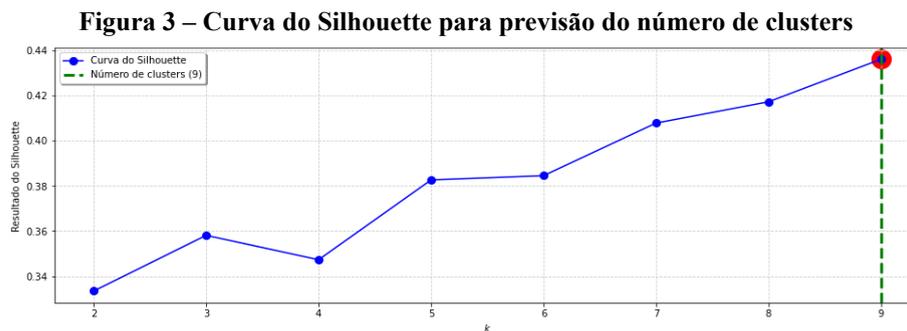
O algoritmo K-Means divide os dados em K clusters, sendo K um valor fornecido pelo usuário (Duda et al., 2001). Dessa forma, para esse algoritmo é necessário definir o número ideal de clusters. Foram utilizados para o problema duas ferramentas com essa finalidade: Elbow Method e Silhouette.

O Elbow Method pode ser impreciso em determinados problemas. Recomenda-se então, utilizar outra métrica de recomendação para número de clusters. Neste trabalho foi utilizado o Silhouette.



Definindo uma distância de -1 a 1, nesse método, o número de clusters ideal é o que mais se aproxima do resultado de 1. O valor é calculado a partir da média do coeficiente Silhouette, representado na Eq. (2), onde a é a distância média das outras distâncias no mesmo cluster. Dessa forma, o número recomendado para o problema foi de 9 clusters, como pode ser visto na Figura 3.

$$s = \frac{b - a}{\max(a, b)} \quad (2)$$



Fonte: autoria própria (2021)

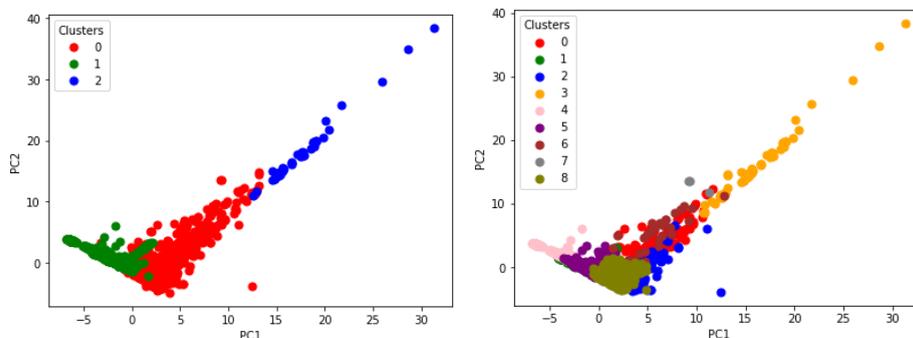
2.3 DBSCAN (do inglês, Density-Based Spatial Clustering of Applications with Noise)

O algoritmo DBSCAN possui como característica sua qualidade para detecção de ruídos, ou *outliers*. É um algoritmo poderoso que descobre clusters de formas e tamanhos arbitrariamente, de acordo com a densidade dos dados presentes, já que clusters são áreas com alta densidade separadas por áreas com baixa densidade (Casas P. et al., 2011).

3 RESULTADOS

Na execução dos códigos, feitos com a linguagem de programação Python na versão 3, a clusterização foi concluída com êxito em ambos os algoritmos de Agrupamento Hierárquico e K-Means, como pode ser visto a partir da Figura 4 em duas e três dimensões.

Figura 4 – Gráfico da ACP com Agrupamento Hierárquico (esquerda) e K-Means (direita)

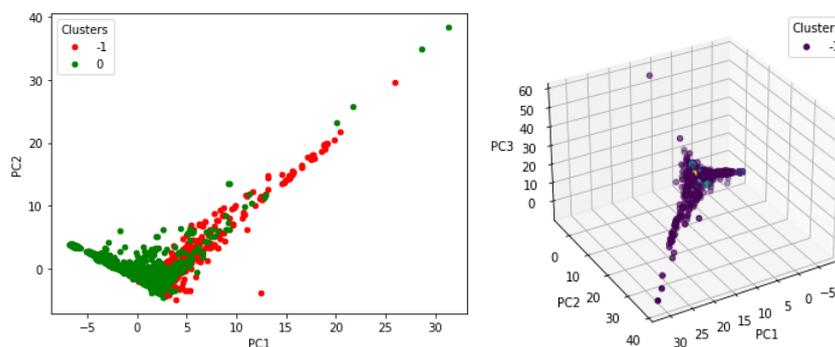




Fonte: autoria própria (2021)

Nota-se uma maior generalização no método de agrupamento hierárquico, pelo número de clusters fornecido, enquanto no K-Means, o alto número de clusters permite uma análise mais específica das características em cada grupo. Já na Figura 5, nota-se uma dificuldade muito grande em detectar os clusters com o DBSCAN nesse problema. Apesar de detectar vários *outliers*, ele não consegue encontrar corretamente os padrões que eram necessários para a formação dos agrupamentos.

Figura 5 – Gráficos da ACP com DBSCAN



Fonte: autoria própria (2021)

A análise final do problema trata da comparação do Quadro 1, contendo alguns exemplos de presença e ausência de ataque na rede, com os resultados obtidos nos gráficos anteriormente. O data frame abaixo possui uma coluna chamada “Label”, que indica ocorrência (1) ou ausência (0) de ataque DDoS.

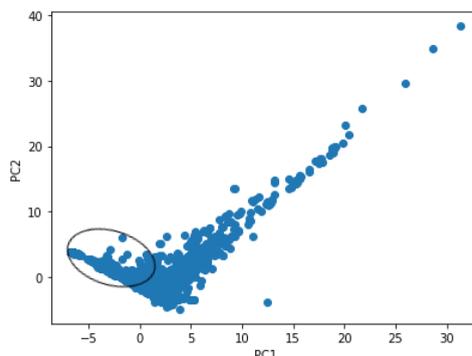
Quadro 1 – Data Frame dos Componentes Principais.

Label	PC1	PC2
1	-1.255697	-0.312235
0	0.545627	-1.550247
1	-3.900416	1.505512
1	-4.597347	2.454929
0	0.213153	-2.069280
0	2.387511	-1.012062

Fonte: autoria própria (2021)

Ao verificar em cada um dos gráficos a coordenada correspondente aos eixos PC1 e PC2, a conclusão obtida é que a maioria dos dados que são ataques estão presentes na região circulada da Figura 6, sendo o cluster 1 da Figura 4 à esquerda e os clusters 1, 5 e 8 da Figura 4 à direita.

Figura 6 – Gráfico da ACP - Região dos ataques



Fonte: autoria própria (2021)

4 CONCLUSÃO

Foi observada a eficácia de alguns algoritmos de aprendizado de máquina não-supervisionados em detectar ataques DDoS. Para o problema em questão e até pela forma que os dados foram apresentados no problema, nota-se que os algoritmos hierárquicos com abordagem aglomerativa e baseados em erro quadrático possuem aplicação com maior facilidade.

Dessa forma, a presente pesquisa pode auxiliar tanto na detecção de anomalias em redes, quanto na escolha de algoritmos para a resolução de problemas semelhantes. Considerando que, neste trabalho, os algoritmos de Agrupamento Hierárquico e K-Means formaram clusters efetivamente, a falta de êxito para o algoritmo baseado em densidade é um item considerado como integrante de trabalhos futuros.

AGRADECIMENTOS

Agradeço à Universidade Tecnológica Federal do Paraná pela oportunidade de participar do Programa Institucional de Voluntariado em Iniciação Científica (PIVIC) e ao professor Luiz Fernando Carvalho pelo apoio, atenção e ensinamentos durante essa pesquisa.

REFERÊNCIAS

- CASAS, P. et al., **UNADA: Unsupervised Network Anomaly Detection Using Sub-space Outliers Ranking**, Springer, Berlin, Heidelberg, Lecture Notes in Computer Science, V. 6640, 2011.
- DUDA, R. et al., **Pattern classification**, New York: John Wiley & Sons, 2001.
- FACELI, Katti; LORENA, Ana Carolina; GAMA, João; CARVALHO, André Carlos Ponce de Leon Ferreira de. **Inteligência artificial: uma abordagem de aprendizado de máquina**. Rio de Janeiro: LTC - Livros Técnicos e Científicos Editora Ltda., 2011.
- PATCHA, Animesh; PARK, Jung-Min. **An overview of anomaly detection techniques: Existing solutions and latest technological trends**, Computer Networks, V. 51, p.3448-3470, 2007.