

Ernane Rosa Martins
(Organizador)

Morris Charts

Line Chart



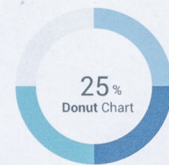
Area Chart



Bar Chart

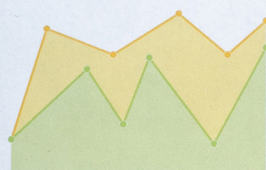


Donut Chart

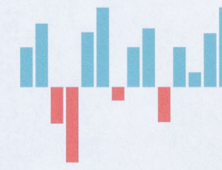


Sparkline Charts

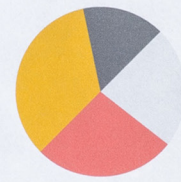
Line Chart



Bar Chart



Pie Chart



Easy Pie Charts



Pesquisa Operacional e sua Atuação Multidisciplinar

Ernane Rosa Martins

(Organizador)

Pesquisa Operacional e sua Atuação Multidisciplinar

**Atena Editora
2019**

2019 by Atena Editora
Copyright © Atena Editora
Copyright do Texto © 2019 Os Autores
Copyright da Edição © 2019 Atena Editora
Editora Executiva: Profª Drª Antonella Carvalho de Oliveira
Diagramação: Karine de Lima
Edição de Arte: Lorena Prestes
Revisão: Os Autores

O conteúdo dos artigos e seus dados em sua forma, correção e confiabilidade são de responsabilidade exclusiva dos autores. Permitido o download da obra e o compartilhamento desde que sejam atribuídos créditos aos autores, mas sem a possibilidade de alterá-la de nenhuma forma ou utilizá-la para fins comerciais.

Conselho Editorial

Ciências Humanas e Sociais Aplicadas

Prof. Dr. Álvaro Augusto de Borba Barreto – Universidade Federal de Pelotas
Prof. Dr. Antonio Carlos Frasson – Universidade Tecnológica Federal do Paraná
Prof. Dr. Antonio Isidro-Filho – Universidade de Brasília
Prof. Dr. Constantino Ribeiro de Oliveira Junior – Universidade Estadual de Ponta Grossa
Profª Drª Cristina Gaio – Universidade de Lisboa
Prof. Dr. Deyvison de Lima Oliveira – Universidade Federal de Rondônia
Prof. Dr. Gilmei Fleck – Universidade Estadual do Oeste do Paraná
Profª Drª Ivone Goulart Lopes – Istituto Internazionele delle Figlie de Maria Ausiliatrice
Prof. Dr. Julio Candido de Meirelles Junior – Universidade Federal Fluminense
Profª Drª Lina Maria Gonçalves – Universidade Federal do Tocantins
Profª Drª Natiéli Piovesan – Instituto Federal do Rio Grande do Norte
Profª Drª Paola Andressa Scortegagna – Universidade Estadual de Ponta Grossa
Prof. Dr. Urandi João Rodrigues Junior – Universidade Federal do Oeste do Pará
Profª Drª Vanessa Bordin Viera – Universidade Federal de Campina Grande
Prof. Dr. Willian Douglas Guilherme – Universidade Federal do Tocantins

Ciências Agrárias e Multidisciplinar

Prof. Dr. Alan Mario Zuffo – Universidade Federal de Mato Grosso do Sul
Prof. Dr. Alexandre Igor Azevedo Pereira – Instituto Federal Goiano
Profª Drª Daiane Garabeli Trojan – Universidade Norte do Paraná
Prof. Dr. Darllan Collins da Cunha e Silva – Universidade Estadual Paulista
Prof. Dr. Fábio Steiner – Universidade Estadual de Mato Grosso do Sul
Profª Drª Girlene Santos de Souza – Universidade Federal do Recôncavo da Bahia
Prof. Dr. Jorge González Aguilera – Universidade Federal de Mato Grosso do Sul
Prof. Dr. Ronilson Freitas de Souza – Universidade do Estado do Pará
Prof. Dr. Valdemar Antonio Paffaro Junior – Universidade Federal de Alfenas

Ciências Biológicas e da Saúde

Prof. Dr. Benedito Rodrigues da Silva Neto – Universidade Federal de Goiás
Prof.ª Dr.ª Elane Schwinden Prudêncio – Universidade Federal de Santa Catarina
Prof. Dr. Gianfábio Pimentel Franco – Universidade Federal de Santa Maria
Prof. Dr. José Max Barbosa de Oliveira Junior – Universidade Federal do Oeste do Pará

Profª Drª Natiéli Piovesan – Instituto Federal do Rio Grande do Norte
Profª Drª Raissa Rachel Salustriano da Silva Matos – Universidade Federal do Maranhão
Profª Drª Vanessa Lima Gonçalves – Universidade Estadual de Ponta Grossa
Profª Drª Vanessa Bordin Viera – Universidade Federal de Campina Grande

Ciências Exatas e da Terra e Engenharias

Prof. Dr. Adélio Alcino Sampaio Castro Machado – Universidade do Porto
Prof. Dr. Eloi Rufato Junior – Universidade Tecnológica Federal do Paraná
Prof. Dr. Fabrício Menezes Ramos – Instituto Federal do Pará
Profª Drª Natiéli Piovesan – Instituto Federal do Rio Grande do Norte
Prof. Dr. Takeshy Tachizawa – Faculdade de Campo Limpo Paulista

Conselho Técnico Científico

Prof. Msc. Abrãao Carvalho Nogueira – Universidade Federal do Espírito Santo
Prof. Dr. Adaylson Wagner Sousa de Vasconcelos – Ordem dos Advogados do Brasil/Seccional Paraíba
Prof. Msc. André Flávio Gonçalves Silva – Universidade Federal do Maranhão
Prof.ª Drª Andreza Lopes – Instituto de Pesquisa e Desenvolvimento Acadêmico
Prof. Msc. Carlos Antônio dos Santos – Universidade Federal Rural do Rio de Janeiro
Prof. Msc. Daniel da Silva Miranda – Universidade Federal do Pará
Prof. Msc. Eliel Constantino da Silva – Universidade Estadual Paulista
Prof.ª Msc. Jaqueline Oliveira Rezende – Universidade Federal de Uberlândia
Prof. Msc. Leonardo Tullio – Universidade Estadual de Ponta Grossa
Prof.ª Msc. Renata Luciane Polsaque Young Blood – UniSecal
Prof. Dr. Welleson Feitosa Gazel – Universidade Paulista

Dados Internacionais de Catalogação na Publicação (CIP) (eDOC BRASIL, Belo Horizonte/MG)	
P474	Pesquisa operacional e sua atuação multidisciplinar [recurso eletrônico] / Organizador Ernane Rosa Martins. – Ponta Grossa, PR: Atena Editora, 2019. Formato: PDF Requisitos de sistema: Adobe Acrobat Reader Modo de acesso: World Wide Web Inclui bibliografia ISBN 978-85-7247-478-8 DOI 10.22533/at.ed.788191107 1. Pesquisa operacional. I. Martins, Ernane Rosa. CDD 658.51
Elaborado por Maurício Amormino Júnior – CRB6/2422	

Atena Editora
Ponta Grossa – Paraná - Brasil
www.atenaeditora.com.br
contato@atenaeditora.com.br

APRESENTAÇÃO

A Pesquisa Operacional (PO) utiliza a matemática, a estatística e a computação para auxiliar na solução de problemas reais, com foco na tomada das melhores decisões nas mais diversas áreas científicas e de atuação humana, buscando otimizar e melhorar suas performances. Através do uso de técnicas de modelagem matemática e eficientes algoritmos computacionais, a PO vem cada vez mais atuando na análise dos mais variados aspectos e situações de problemas complexos em demandas de inúmeras áreas, principalmente por conta de sua flexibilidade de aplicação e interação multidisciplinar, permitindo a tomada de decisões efetivas e a construção de sistemas mais produtivos.

Esta obra reúne importantes trabalhos que envolvem o uso de PO, realizados em diversas instituições de ensino do Brasil, abordando assuntos atuais e relevantes, tais como: modelos matemáticos; otimização multiobjectivo; heurísticas; algoritmos; otimização geométrica; metodologia SODA; soft systems methodology; strategic choice approach; procedimentos metodológicos de análise estatística; jogos cooperativos; algoritmos genéticos; método VIKOR; regressão linear múltipla; algoritmos de aprendizado de máquina; análise de decisão multicritério e composição probabilística de preferências.

A importância desta coletânea está na excelência dos trabalhos apresentados e na contribuição dos seus autores em temas de experiências e vivências. A socialização destes estudos no meio acadêmico, permite ampla análise e inúmeras discussões sobre diversos assuntos pertinentes referentes a atuação multidisciplinar da PO. Por fim, agradeço a todos que contribuíram na construção desta belíssima obra e desejo a todos os leitores, boas reflexões sobre os assuntos abordados.

Ernane Rosa Martins

SUMÁRIO

CAPÍTULO 1	1
UMA ABORDAGEM MULTIOBJETIVO EM UM PROBLEMA DE PRODUÇÃO COM ESTOQUE INTERMEDIÁRIO E TESTE DE FUNCIONALIDADE	
Sander Joner Neida Maria Patias Volpi Joyce Rodrigues da Silva Tulipa Gabriela Guilhermina Juvenal da Silva	
DOI 10.22533/at.ed.7881911071	
CAPÍTULO 2	16
SOLUÇÕES INTEIRAS PARA O PROBLEMA DE CORTE DE ESTOQUE UNIDIMENSIONAL	
Gonçalo Renildo Lima Cerqueira Sérgio da Silva Aguiar Marlos Marques	
DOI 10.22533/at.ed.7881911072	
CAPÍTULO 3	28
OTIMIZAÇÃO GEOMÉTRICA DE AERONAVES REMOTAMENTE PILOTADAS CARGUEIRAS VIA ECOLOCALIZAÇÃO	
Guilherme Aparecido Barbosa Pereira Ivo Chaves da Silva Júnior Luiz Rogério Andrade de Oliveira Carlos Henrique Sant'Ana da Silva	
DOI 10.22533/at.ed.7881911073	
CAPÍTULO 4	41
O CASO DA INDÚSTRIA CRIATIVA DO CARNAVAL SOB O ENFOQUE DO SODA	
Ailson Renan Santos Picanço Adjame Alexandre Oliveira Mischel C.N. Belderrain Nissia Carvalho Rosa Bergiante	
DOI 10.22533/at.ed.7881911074	
CAPÍTULO 5	55
MODELO DE NEGÓCIO EM UMA COMUNIDADE AGRÍCOLA: APLICAÇÃO DE <i>SOFT SYSTEMS METHODOLOGY</i> E <i>STRATEGIC CHOICE APPROACH</i>	
Michelle Carvalho Galvão Silva Pinto Bandeira Raquel Issa Mattos Mischel Carmen Neyra Belderrain Anderson Ribeiro Correia John Bernhard Kleba	
DOI 10.22533/at.ed.7881911075	
CAPÍTULO 6	72
MODELAGEM MATEMÁTICA PARA GERAÇÃO DE ESCALAS DE TURNO	
Laiz de Carvalho Nogueira Tiago Araújo Neves	
DOI 10.22533/at.ed.7881911076	

CAPÍTULO 7	87
METODOLOGIA ADOTADA PELA ARCELORMITTAL BRASIL PARA CERTIFICAÇÃO DE PADRÕES SECUNDÁRIOS PARA ANÁLISES QUÍMICAS EM AMOSTRAS DE MINÉRIO DE FERRO DA MINA DE SERRA AZUL EM MINAS GERAIS	
Antonio Fernando Pêgo e Silva Juliana Cecília C R Vieira Luiz Paulo de Carvalho Serrano	
DOI 10.22533/at.ed.7881911077	
CAPÍTULO 8	100
JOGOS COOPERATIVOS NA ALOCAÇÃO DE CUSTOS DE ESTOQUES DE PEÇAS COMPARTILHADOS	
Bernardo Santos Aflalo Natália Nogueira Ferreira Souza Takashi Yoneyama	
DOI 10.22533/at.ed.7881911078	
CAPÍTULO 9	112
BIASED RANDOM-KEY GENETIC ALGORITHM ACCORDING TO LEVY DISTRIBUTION FOR GLOBAL OPTIMIZATION	
Mariana Alves Moura Ricardo Martins de Abreu Silva	
DOI 10.22533/at.ed.7881911079	
CAPÍTULO 10	126
AVALIAÇÃO MULTICRITÉRIO DA QUALIDADE DA INFORMAÇÃO CONTÁBIL	
Alini da Silva Nelson Hein Adriana Kroenke	
DOI 10.22533/at.ed.78819110710	
CAPÍTULO 11	142
AVALIAÇÃO DE MODELOS COMPUTACIONAIS DE APRENDIZADO DE MÁQUINA PARA DETECÇÃO REATIVA E PREVENTIVA DE BOTNETS	
Vinicius Oliveira de Souza Sidney Cunha de Lucena	
DOI 10.22533/at.ed.78819110711	
CAPÍTULO 12	158
AVALIAÇÃO DE ATRIBUTOS ESTATÍSTICOS NA DETECÇÃO DE ATAQUES DDOS BASEADA EM APRENDIZADO DE MÁQUINA	
Eduardo da Costa da Silva Sidney Cunha de Lucena	
DOI 10.22533/at.ed.78819110712	

CAPÍTULO 13	173
ABORDAGEM PROBABILÍSTICA À ESCOLHA DE PRODUTOS DE DEFESA: UMA APLICAÇÃO DA COMPOSIÇÃO PROBABILÍSTICA DE PREFERÊNCIAS NA AQUISIÇÃO DE BLINDADOS	
Luiz Octávio Gavião	
Annibal Parracho Sant'Anna	
Gilson Brito Alves Lima	
Pauli Adriano de Almada Garcia	
DOI 10.22533/at.ed.78819110713	
CAPÍTULO 14	189
A STOCHASTIC DYNAMIC MODEL FOR SUPPORT OF THE MANAGEMENT OF NEW PRODUCT DEVELOPMENT PORTFOLIOS	
Samuel Martins Drei	
Thiago Augusto de Oliveira Silva	
Marco Antonio Bonelli Júnior	
Luciana Paula Reis	
Matheus Correia Teixeira	
DOI 10.22533/at.ed.78819110714	
CAPÍTULO 15	205
A RELAXED FLOW-BASED FORMULATION FOR THE OPEN CAPACITATED ARC ROUTING PROBLEM	
Rafael Kendy Arakaki	
Fábio Luiz Usberti	
DOI 10.22533/at.ed.78819110715	
CAPÍTULO 16	217
A COMPOSIÇÃO PROBABILÍSTICA DE PREFERÊNCIAS COM MEDIDAS DE DESIGUALDADE: CORRELAÇÕES COM OS PONTOS DE VISTA PROGRESSISTA E CONSERVADOR	
Luiz Octávio Gavião	
Annibal Parracho Sant'Anna	
Gilson Brito Alves Lima	
DOI 10.22533/at.ed.78819110716	
SOBRE O ORGANIZADOR	233

AVALIAÇÃO DE ATRIBUTOS ESTATÍSTICOS NA DETECÇÃO DE ATAQUES DDOS BASEADA EM APRENDIZADO DE MÁQUINA

Eduardo da Costa da Silva

Universidade Federal do Estado do Rio de Janeiro (UNIRIO), Programa de Pós-Graduação de Informática

Rio de Janeiro –RJ

50° Simpósio Brasileiro de Pesquisa Operacional (SBPO 2018), ISSN 1518-1731

Sidney Cunha de Lucena

Universidade Federal do Estado do Rio de Janeiro (UNIRIO), Programa de Pós-Graduação de Informática

Rio de Janeiro - RJ

50° Simpósio Brasileiro de Pesquisa Operacional (SBPO 2018), ISSN 1518-1731

RESUMO: Identificar ataques a serviços de TI é tarefa difícil, dada a grande quantidade de fluxos de dados numa rede. Ademais, há diversos casos onde o tráfego de ataque tem padrão similar ao tráfego legítimo. Logo, modelar eficientemente as características do tráfego da rede se faz necessário para destacar anomalias que possam indicar um ataque. Este trabalho analisa o impacto que determinados atributos do tráfego em uma rede exercem no desempenho de alguns sistemas de detecção de ataques baseados em aprendizado de máquina. Em especial, propõe-se o uso de atributos obtidos a partir de medidas estatísticas envolvendo o agregado de fluxos de dados coexistentes

num intervalo de tempo. Para esta análise, foi implementado um sistema de detecção baseado na arquitetura *lambda* e os experimentos usaram *datasets* realísticos. Os resultados demonstram que a inclusão dos atributos propostos pode melhorar o desempenho de certos algoritmos de aprendizado de máquina na detecção de ataques.

PALAVRAS-CHAVE: Detecção de anomalia, aprendizado de máquina, entropia.

ABSTRACT: Identifying attacks towards IT services can be very difficult due to the large amount of data flow in a network. In addition, there are several cases where malicious and legitimate traffics have similar patterns. Therefore, it is necessary to efficiently model the network traffic in order to highlight any anomaly that may indicate an attack. This work analyzes the impact of certain network traffic attributes on the performance of some attack detection systems based on machine learning. In particular, it proposes the use of attributes obtained from statistical measures involving the aggregated data flows that coexist in a given time interval. The analysis was performed through the implementation of a detection system based on the *lambda* architecture and experiments used realistic datasets. The results show that the inclusion of the proposed attributes may increase the performance of certain machine

learning algorithms for attack detection.

KEYWORDS: Anomaly detection, machine learning, entropy.

1 | INTRODUÇÃO

Durante os últimos anos, o tráfego de dados na Internet tem apresentado um acentuado crescimento. Em 2010, por exemplo, foram gerados 1 bilhão de *gigabytes* de dados a cada 2 dias e espera-se que até 2020 sejam criados e manipulados mais de 40 ZB [Yi et al., 2014]. Paralelamente, vem aumentando também a quantidade de tráfego ilícito gerado por atividades maliciosas. É o caso, dentro outros, dos ataques de negação de serviço distribuídos (DDoS - *Distributed Denial of Service*), que visam esgotar os recursos computacionais da vítima ou da rede que a conecta a partir de fluxos de pacotes de dados com requisições forjadas e origens distribuídas pela Internet. Logo, garantir que uma rede de computadores esteja operando satisfatoriamente, livre de tráfegos maliciosos oriundos de ataques, tornou-se uma atividade primordial. No entanto, há dificuldades no que tange à análise do tráfego de rede para a construção de padrões comportamentais que permitam uma correta identificação dos ataques, que podem ser caracterizados como anomalias no tráfego. Tais dificuldades derivam não apenas da grande quantidade de fluxos de dados e suas peculiares, mas também do avanço das técnicas empregadas nos ataques, que objetivam camuflar o tráfego ilícito para dificultar sua identificação [Ficco and Rak, 2015]. Deste modo, esta modelagem deve considerar atributos e métricas, relacionadas aos fluxos de dados que passam por uma rede, que sejam relevantes ao desempenho dos mecanismos empregados na análise do tráfego.

Diante desse contexto, este artigo investiga certos mecanismos de detecção de ataques baseados em em árvore de decisão, no modo *on-line*, para analisar o desempenho destes sistemas com relação a certos atributos associados aos fluxos de dados que passam pela rede que se deseja proteger. Em especial, objetiva-se verificar se é possível melhorar este desempenho a partir do acréscimo de atributos relacionados a medidas estatísticas extraídas para o tráfego agregado da rede num dado intervalo de tempo. Para tal, foi proposto e implementado um sistema baseado na arquitetura *lambda* [Marz and Warren, 2015], que considera o atual cenário tráfego de dados e alinha-se com o conceito de *Big Data* para gerar respostas em tempo hábil. A avaliação do sistema utilizou um *dataset* realístico formado por uma mescla de tráfego real legítimo, ataques DDoS gerados a partir de um ambiente controlado e dados relativos a três diferentes cenários do *dataset* CTU-13 [Garcia et al., 2014]. Os resultados obtidos demonstram que a inclusão dos atributos propostos pode melhorar o desempenho de certos algoritmos de aprendizado de máquina na detecção de ataques DDoS.

A estrutura deste artigo é a seguinte: a Seção 2 apresenta alguns conceitos básicos; a Seção 3 apresenta trabalhos relacionados; na Seção 4 tem-se arquitetura

do sistema desenvolvido; a avaliação experimental é descrita na Seção 5; a Seção 6 traz os resultados obtidos; e a Seção 7 conclui artigo.

2 | CONCEITOS BÁSICOS

Esta seção descreve como se dá uma classificação de dados utilizando aprendizado de máquina baseado em árvore de decisão e as características da entropia de *Shannon*, medida estatística usada em um dos atributos avaliados.

2.1 Árvores de Decisão

Um modelo de classificação que utiliza a estrutura de dados árvore de decisão é determinado como um processo de classificação que faz uma divisão recursiva de um conjunto de dados em conjuntos menores, divisão esta baseada em regras simples a partir dos valores dos dados. Essas divisões ocorrem nos nós de uma estrutura em árvore, com as consequentes ramificações podendo ou não levar a outros nós até se atingir o subconjunto mínimo, representado como uma *folha*. As regras de divisão contidas nos nós são automaticamente geradas a partir de uma etapa de treinamento, que necessariamente precisa contar com um *dataset* marcado, ou seja, cujo conteúdo é completamente descrito pela inserção de *labels* em suas instâncias que informam qual categoria a instância de fato representa. Tal estratégia de aprendizado de máquina, com a necessidade de marcações, é chamada de supervisionada.

Dentre os tipos de mecanismos de aprendizado supervisionado de máquina existentes, como por exemplo, Redes Neurais e Máquina de Vetores de Suporte, optou-se pelo uso de árvores de decisão por sua simplicidade e eficiência em capturar padrões complexos [Michie et al., 1994], principalmente quando se considera a necessidade destes mecanismos trabalharem em tempo real.

Os mecanismos de árvore de decisão usados neste trabalho foram o VHT (*Vertical Hoeffding Tree*), o *Bagging* e o *Adaptive Bagging* [Quinlan, 1986], todos capazes de efetuar aprendizado de máquina com base num fluxo de dados (modo *on-line*) ou incremental. O VHT é um classificador distribuído e escalável que utiliza paralelismo vertical baseado em uma árvore de Hoeffding ou *Very Fast Decision Tree* (VFDT). O método *Bagging* busca maximizar a acurácia pelo agrupamento de diferentes árvores de decisão, superando erros de classificação ao se usar apenas um classificador [Khan et al., 2007]. O modelo final de classificação é resultado de uma combinação dos resultados preditivos individuais desses outros métodos. Já o *Adaptive Bagging* tem funcionamento análogo ao *Bagging*, exceto pelo fato de monitorar a acurácia da classificação ao longo da fase de treinamento e de refinar o seu modelo com os últimos dados capturados. É indicado para situações não-estacionárias de aprendizado, onde os padrões dos fluxos de dados podem variar com o decorrer do tempo, como pode ocorrer com os fluxos de dados numa rede.

2.2 Entropia de Shannon

A entropia de Shannon, enquanto métrica estatística, é usada neste trabalho para descrever o grau de dispersão ou concentração das distribuições de probabilidades de certas características de rede. Essa métrica é bastante adotada na literatura para a identificação de tráfego malicioso, propiciando uma visão de estado global das características de rede dentro de um intervalo de tempo. É definida pela seguinte equação [Shannon, 1948]:

$$E = - \sum_{i=1}^N p_i \log_2 (p_i), \quad (1)$$

onde N é a quantidade de diferentes ocorrências i no espaço amostral e p_i é a probabilidade associada a cada ocorrência i .

3 | TRABALHOS RELACIONADOS

O trabalho de [Bhuyan et al., 2014] traz um estudo comparativo entre várias técnicas para detecção de anomalias, ressaltando a efetividade do uso de aprendizado de máquina para este fim. Seguindo a mesma tendência de utilização de aprendizado de máquina, [Singh et al., 2014] faz a detecção de anomalias oriundas de uma *botnets*. O Apache Hadoop, o Apache Hive e o Apache Mahout compõem os sistemas de uma arquitetura projetada para ser escalável e lidar com grande quantidade de dados. Verificou-se que, embora seja possível analisar fluxos de dados com o Hadoop, a latência provocada pelo processamento em lote comprometia a velocidade de execução do sistema.

Já utilizando o processamento em tempo real para analisar anomalias, geradas por ataques, no momento de suas ocorrências, os autores em [Du et al., 2014] usam o Apache Storm [Toshniwal et al., 2014]. Os resultados indicam que a proposta apresenta bom desempenho e escalabilidade.

Em [Robinson and Thomas, 2015], foram comparados dez algoritmos supervisionados de aprendizado de máquina, no modo *batch*, para a classificação de ataques. A avaliação usou três *datasets* antigos, já conhecidos: LLDDoS, CAIDA DDoS 2007 e CAIDA Conficker. O algoritmo Adaboost, com classificador base *Random Forest*, obteve o melhor desempenho, alcançando acurácia máxima de 99,96% no *dataset* CAIDA. Mais uma vez, em [Lobato et al., 2016] ataques DoS e de varredura são detectados através de classificação supervisionada usando os algoritmos C4.5, rede neural e máquina de suporte de vetores. Anomalias são sinalizadas assumindo-se uma distribuição Gaussiana nos valores dos atributos mais relevantes e um limiar empírico. Apesar do sistema proposto alcançar mais de 95% de acurácia, a classificação dos fluxos em normais, ataques DoS e de varredura não leva em consideração o estado global do tráfego da rede.

Conforme observado nos trabalhos supracitados, ao longo dos últimos anos diversas técnicas de aprendizado de máquina vêm sendo pesquisadas para fins de detecção de ataques. Porém, diferentemente desses estudos, o presente artigo visa destacar como determinados atributos estatísticos obtidos dos fluxos de dados podem influenciar certos métodos de aprendizado de máquina *on-line* utilizados para a detecção de ataques DDoS em volumosos fluxos de dados. Além disso, foram também levados em consideração aspectos operacionais relacionados com o sistema de detecção, tais como tempo disponível para análise dos fluxos, volume de tráfego e variedade de comportamento ao longo do tempo, tudo isto representado por *datasets* atuais e realísticos.

4 | ARQUITETURA DO SISTEMA DE DETECÇÃO DE ATAQUES

A arquitetura do sistema proposto é baseado na arquitetura *Lambda*, cujo arranjo objetiva processar, em tempo hábil, volumosos fluxos de dados. Sua estrutura é composta de 3 camadas: processamento em velocidade, processamento em lote e serviços. A camada de processamento em lote é responsável por armazenar grande quantidade de dados e processá-los em lote. A camada de processamento em velocidade é designada para a manipulação de fluxos de dados em alta velocidade. Já a camada de serviço tem a função de juntar os resultados dos processamentos das outras duas camadas e disponibilizar os dados analisados. Em adição, foi desenvolvido um módulo denominado “holístico”, cuja principal função é adicionar características temporais do estado global da rede a cada intervalo de tempo. Além disso, esse módulo também computa as métricas de desempenho da classificação.

A camada de processamento em alta velocidade é composta pelas ferramentas Apache SAMOA [De Francisci Morales, 2013] e Apache Storm [Toshniwal et al., 2014], que por sua vez usam a camada de processamento em lote, representada pelo Apache Hadoop [Manikandan and Ravi, 2014], para obtenção de fluxos de treinamento e teste. Funções internas do módulo holístico fazem o papel da camada de serviço, resumizando os dados das duas outras camadas para posterior exibição das métricas de desempenho. A Figura 1, abaixo, exhibe a arquitetura do sistema.

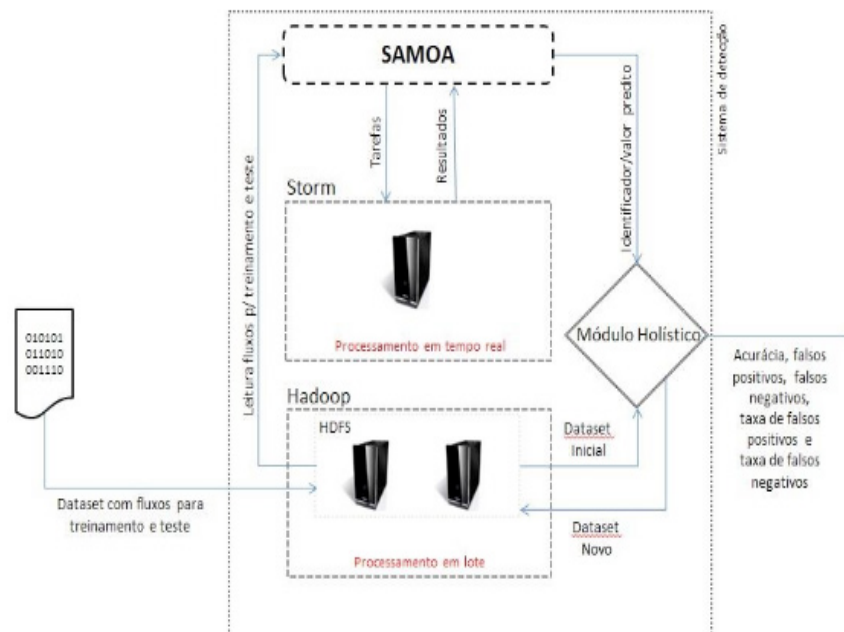


Figura 1. Arquitetura do sistema

4.1 Componentes e Funcionamento do Sistema

O Apache SAMOA é uma plataforma destinada a encontrar padrões em grandes volumes de fluxos de dados, fornecendo algoritmos para as mais comuns tarefas de mineração de dados e aprendizado de máquina. Já o Apache Storm é um sistema distribuído para processamento em tempo real de volumosos fluxos de dados, enquanto que o Apache Hadoop é destinado a armazenamento e processamento em lote, sendo que somente é utilizado seu módulo HDFS, que serve como repositório para os dados manipulados pelo sistema.

O módulo holístico tem a responsabilidade de ler os fluxos do *dataset* armazenado no HDFS e dele extrair novos atributos com características temporais, criando assim um novo *dataset* enriquecido com tais informações. Estas características temporais estão diretamente associadas ao estado global dos fluxos da rede a cada 30 segundos, que é um intervalo de tempo adotado em trabalhos similares [Joldzic et al., 2016]. Estes estados globais são representados pelas entropias associadas aos endereços IP de origem, números das portas de protocolo de origem, endereços IP de destino e numeração das portas de destino. Também são registradas as suas respectivas diferenças para os valores de entropias calculados nos quatro intervalos de tempo anteriores, registrando-se assim a variação dos valores das entropias ao longo dos últimos dois minutos [Joldzic et al., 2016]. Além das entropias, calculam-se também o tamanho médio dos fluxos em *bytes*, o número médio de pacotes dos fluxos e o desvio-padrão dos tamanhos dos pacotes e das durações dos fluxos existentes a cada 30 segundos. Portanto, a cada fluxo são acrescentados esses 20 atributos contendo características estatísticas temporais relativas ao agregado de tráfego fluindo no respectivo intervalo de 30 segundos ao qual o fluxo pertence.

O funcionamento do sistema inicia-se pelo armazenamento do dataset no

Hadoop, que passa a ser chamado de dataset inicial. O módulo holístico então lê os fluxos do dataset inicial, agrupando-os em intervalos de tempo de 30 segundos para cálculo dos atributos estatísticos (denominados de atributos avançados). Posteriormente, o módulo holístico acrescenta os novos atributos no dataset e o disponibiliza ao SAMOA. Com base nesse novo dataset, o SAMOA cria uma tarefa a ser enviada ao STORM para fins de processamento dos fluxos. A medida que os fluxos vão sendo analisados, os resultados das classificações - predição em “normal” ou “ataque” e identificador do fluxo são enviados para o módulo holístico. O módulo holístico, alicerçado em marcações anteriores ao processo de classificação, que determinam se o fluxo é de fato ataque ou não, calcula as métricas de desempenho do sistema de detecção e disponibiliza em arquivo.

O sistema foi implementado em um conjunto de três máquinas virtuais, uma com 20GB de RAM e as outras duas com 6GB de RAM cada. Todas as máquinas utilizaram sistema operacional Debian 7 e quatro núcleos de processador Xeon(R) CPU E5-2643 v2 @ 3.50GHz.

5 | AVALIAÇÃO EXPERIMENTAL

Nesta seção serão descritos os procedimentos para a construção dos *datasets*, a organização dos experimentos e as métricas analisadas.

5.1 Criação dos Datasets

Os *datasets* criados para avaliar o sistema proposto derivam de um *dataset* realístico, obtido a partir da combinação de tráfego real e tráfego de ataque emulado, além de um conjunto de *datasets* publicamente disponível conhecido como CTU-13.

O CTU-13 [Garcia et al., 2014] compreende um conjunto de 13 diferentes cenários contendo *botnets* conhecidas de variados tipos. Num primeiro momento, os *datasets* gerados para o conjunto CTU-13 continham capturas de fluxos unidirecionais [Garcia et al., 2014]. Posteriormente, os autores geraram outro conjunto de capturas para os mesmos cenários, mas desta vez para fluxos bidirecionais [Garcia and Uhler, 2017]. No presente trabalho, tanto no conjunto de fluxos unidirecionais quanto no conjunto de fluxos bidirecionais foram usados os cenários 4, 10 e 11, por conterem ataques DDoS.

O *dataset* realístico foi gerado combinando dados de tráfego legítimo (sem ataques), obtidos de uma rede real, e dados de tráfego de ataque DDoS obtidos por emulação. O tráfego legítimo foi obtido da rede de computadores do Observatório Nacional (ON) durante 7 dias. A comprovação da ausência de ataques baseia-se nas informações fornecidas pelo sistema de segurança de rede do ON. O tráfego de ataque DDoS foi gerado através da ferramenta *hping3* em uma rede isolada com duas máquinas. A origem distribuída dos ataques foi emulada forjando-se e amplamente variando os IPs de origem dos pacotes gerados pelo *hping3*. Os ataques emulados

representaram dois tipos de ataques DDoS de alta incidência: tráfego UDP para a porta de destino 53 e tráfego ICMP.

Vale salientar que a literatura não define proporções exatas entre a quantidade de fluxos de tráfego normal e a quantidade de fluxos de ataque verificados no momento de um ataque. Para o *dataset* realístico, foi escolhida a proporção de 80% de fluxos normais e 20% de fluxos de ataques para a mesclagem desses arquivos. Essa proporção visa dificultar a obtenção, via mecanismos de aprendizados de máquina, dos padrões ocultos de comportamento dos fluxos em virtude do desbalanceamento entre a quantidade de fluxos de treinamento e a quantidade de fluxos de teste.

A partir dos *datasets* citados, foram criados três novos grupos de *datasets* contendo fluxos de treinamento e fluxos de teste, onde cada grupo possui um conjunto separado para diferentes cenários do CTU-13: Grupo 1 - tanto os fluxos de treinamento quanto os de teste são oriundos dos *datasets* CTU-13; Grupo 2 - fluxos de treinamento provenientes do *dataset* realístico e fluxos de teste provenientes do *dataset* CTU-13 com fluxos unidirecionais; e Grupo 3 - fluxos de treinamento provenientes do *dataset* realístico e fluxos de teste provenientes do *dataset* CTU-13 com fluxos bidirecionais.

Esta divisão em grupos busca representar situações diferentes. O Grupo 1 reflete uma prática costumeira em estudos de classificadores usando aprendizado de máquina, onde tanto os dados de treinamento quanto os de teste se originam de diferentes partes de um mesmo *dataset*. Já os Grupos 2 e 3 buscam representar uma situação de cunho prático e operacional na administração de uma rede, situação esta onde a disponibilidade de *datasets* marcados para fins de treinamento só é possível a partir de fontes externas à instituição. Trata-se de uma situação análoga às assinaturas de tráfego contendo ataques recém-descobertos que precisam ser inseridas nos sistemas de detecção de intrusão.

5.2 Organização dos Experimentos

Cada um dos três grupos de *datasets* gerou quatro diferentes conjuntos, cada qual contendo um diferente conjunto de atributos associados aos fluxos IP: (i) um conjunto chamado de *dataset* puro ou primário, contendo apenas os atributos de fluxo gerados pelo *nProbe*, aqui chamados de atributos primários; (ii) um conjunto chamado de *dataset* com estatísticas, contendo os atributos primários e incluindo o tamanho médio dos fluxos em bytes, tamanho médio dos fluxos em pacotes, desvio-padrão das durações dos fluxos e desvio-padrão dos tamanhos dos pacotes dos fluxos, todos calculados para o tráfego agregado da rede; (iii) um conjunto chamado de *dataset* com entropias, contendo os atributos primários e incluindo as medidas de entropia do tráfego agregado e suas diferenças em relação aos quatro intervalos anteriores de 30 segundos; por fim, (iv) um conjunto chamado de *dataset* completo, contendo todos os atributos gerados. Os atributos primários, gerados pelo *nProbe*, são os seguintes: marca de tempo de início, duração do fluxo, protocolo (campo do cabeçalho IP), IP de

origem, porta de origem, IP de destino, porta de destino, número de bytes e número de pacotes do fluxo.

Uma vez obtidos esses quatro conjuntos de atributos para cada um dos três grupos de *datasets*, foram aplicados, em cada conjunto, os métodos de classificação supervisionada VHT, *Bagging* e *Adaptive Bagging* - todos oferecidos pelo SAMOA. Considerando o processo de classificação de fluxos de dados cujo o fator tempo é essencial, foi utilizada a abordagem *interleaved-test-then-train* como técnica de validação da classificação

- uma das mais utilizadas em cenários de classificação de fluxos de dados por não onerar demasiadamente o tempo de execução, ao contrário do que acontece, por exemplo, com a validação cruzada.

5.3 Métricas de Desempenho Utilizadas

As métricas utilizadas neste trabalho são aquelas normalmente encontradas na literatura para a aferição da eficiência de um sistema de detecção de ataques [Bhuyan et al., 2014]. Todas elas se baseiam nas seguintes variáveis: Falsos positivos (FP) - quantidade de fluxos normais classificados como fluxos de ataque; Falsos negativos (FN) - quantidade de fluxos de ataque classificados como fluxos normais; Verdadeiros positivos (VP) - quantidade de fluxos de ataque classificados como fluxos de ataques; e Verdadeiros negativos

(VN) - quantidade de fluxos de normais classificados como fluxos normais. A partir dessas medidas, tem-se as seguintes métricas:

- Acurácia - razão percentual entre a quantidade de fluxos corretamente classificados e a quantidade total de fluxos analisados:

$$\text{Acurácia} = \frac{\text{qtd de fluxos analisados} - (\text{FP} + \text{FN})}{2\text{qtd total de fluxos analisados}} * 100 \quad (2)$$

- Precisão - razão percentual entre a quantidade de fluxos de ataque corretamente classificados (VP) e a quantidade de fluxos classificados como ataque (VP + FP):

$$\text{Precisão} = \frac{\text{VP}}{\text{VP} + \text{FP}} * 100 \quad (3)$$

- Taxa de verdadeiros positivos, ou *Recall* - razão percentual entre o número de verdadeiros positivos (VP) e a quantidade de fluxos de ataque:

$$\text{Recall} = \frac{\text{VP}}{\text{qtd total de fluxos de ataque}} * 100 \quad (4)$$

- *F-Measure* - média harmônica entre precisão e *recall*, que, neste trabalho, assume valores entre 0 e 10000, onde o primeiro valor representa o pior resultado e o segundo, o melhor:

$$F\text{-Measure} = \frac{2}{\frac{1}{\text{Precisão}} + \frac{1}{\text{Recall}}} * 100 \quad (5)$$

6 | RESULTADOS

Nesta seção serão apresentados os impactos nas métricas de desempenho elencadas, considerando a variação dos arranjos de atributos nos Grupos 1, 2 e 3. Conforme já mencionado, o VHT, o *Bagging* e o *Adaptive Bagging* foram os mecanismos de aprendizado de máquina utilizados no processo de classificação ou identificação dos fluxos como “normal” ou “ataque”. A variação desses conjuntos de atributos tem a intenção representar as diferentes contribuições das medidas estatísticas provenientes do módulo holístico na detecção dos ataques.

Nos Grupo 1 e Grupo 2, nas quatro variações dos conjuntos de atributos dos cenários 4, 10 e 11 e utilizando os três mecanismos de aprendizado de máquina, foram obtidos 100% de acurácia, 100% de precisão, 100% de taxa de verdadeiros positivos, ou *recall*, 100% de taxa de verdadeiros negativos e valor máximo para a *F-Measure*, não havendo classificação errada de nenhum fluxo normal como fluxo de ataque e vice-versa.

Já para o Grupo 3, que contempla fluxos bidirecionais e que usa *datasets* de fontes diferentes para treinamento e teste, ao se usar o mecanismo VHT os resultados mostram que, pela primeira vez, o acerto máximo não foi alcançado. Verificou-se 99,84% de acurácia para todas as variações no cenário 4, 94,32% de acurácia para todas as variações no cenário 10, e 94,67% de acurácia para todas as variações no cenário 11. Tanto a acurácia, a precisão, o *F-Measure* e a taxa de falsos positivos foram influenciados pela quantidade de fluxos normais classificados erroneamente como fluxos de ataques. A maior influência observada foi na precisão, pois essa medida considera todos os casos nos quais os fluxos são classificados como positivos para ataque, sendo ou não verdadeiramente positivos, tendo como resultado 59,16% para o cenário 4 e de 58,83% para os cenários 10 e 11. A precisão, por sua vez, interferiu na *F-Measure*, com 7434,09 para todas as variações no cenário 4, 7408,06 para todas as variações no cenário 10 e 7408,01 para todas as variações no cenário 11. No entanto, não houve casos de fluxos de ataques classificados erroneamente como fluxos normais, o equivalente a falsos negativos. Por conseguinte, todos os fluxos de ataques foram corretamente identificados, resultando numa taxa de verdadeiros positivos, ou *recall*, igual a 100%. Ainda como consequência da quantidade de falsos positivos, foi obtida uma taxa de verdadeiros negativos igual a 99,84% para o cenário 4, 93,82%

para o cenário 10 e 94,23% para o cenário 11. Ainda com relação aos resultados do VHT para o Grupo 3, verificou-se que a variação da natureza dos atributos não interferiu nos resultados. Qualquer que fosse a combinação de atributos dentro do mesmo cenário, entre os primários e os gerados pelo módulo holístico, não houve variação nos valores das métricas empregadas.

Tal qual ocorrido para o VHT, no método *Bagging* os resultados não se modificaram com a inclusão dos atributos gerados pelo módulo holístico, a exceção do cenário 11. Neste cenário, é possível verificar uma melhora em todas as métricas de desempenho quando é usado o conjunto completo de atributos. Foi verificado, também, que a acurácia passou de 94,67% para 100%, a precisão passou de 58,83% para 99,96%, a *F-Measure* foi de 7408,01 para 9998,16, a taxa de verdadeiros negativos de 94,23% para 100%, a quantidade de falsos positivos diminuiu de 5.713 para 3, a quantidade de verdadeiros negativos sofreu um aumento de 5.710 casos e a taxa de falsos positivos apresentou uma queda de 5,77% para 0,003%.

Assim como ocorreu nos Grupos 1 e 2, quando o mecanismo *Adaptive Bagging* foi aplicado ao Grupo 3, independente do cenário analisado e das variações dos conjuntos de atributos, obteve-se eficiência máxima nas métricas de desempenho.

Os valores ótimos das métricas de desempenho obtidos para o Grupo 1 podem ser justificados pela utilização de partes de um mesmo *dataset*, gerado num mesmo ambiente de rede, para treinamento e teste do sistema, levando a crer que os fluxos IPs acabam por apresentar um padrão mais definido de comportamento. Portanto, pode-se concluir que, neste caso, as variações dos conjuntos de atributos tornam-se sem efeito para os mecanismos de aprendizado de máquina, que conseguem capturar mais facilmente esses padrões comportamentais.

No caso do Grupo 2, apesar de se considerar *datasets* diferentes para treinamento e teste, o fato do *dataset* de teste ser proveniente do CTU-13 com fluxos IPs unidirecionais pode justificar o acerto total das classificações apenas usando-se os atributos primários. Os pesquisadores responsáveis pelo CTU-13 alegam em [Garcia and Uhler, 2017] que o *dataset* contendo dados bidirecionais possui informações mais representativas a respeito do tráfego de uma rede real e explicitamente recomendam que este *dataset* seja usado no lugar do *dataset* unidirecional. Pode-se então concluir que o *dataset* unidirecional apresenta um padrão de comportamento mais simples, o que possivelmente facilita a classificação. Apesar disso, os resultados relativos ao *dataset* unidirecional foram aqui mantidos para uma maior entendimento a respeito da eficiência do sistema proposto.

No Grupo 3, onde os experimentos usaram fluxos de teste provenientes do *dataset* CTU-13 bidirecional, os dados armazenados tendem a refletir melhor a diversidade e a dinâmica dos fluxos, resultando numa maior variação do padrão comportamental dos atributos primários. Consequentemente, a tarefa de classificação passa a ser mais complexa, ainda mais quando se considera que os dados de treinamento são originados de um *dataset* obtido de outro ambiente de rede.

Tal característica, somada ao fato da haver grande diferença (ou desbalanceamento) na quantidade de fluxos entre as duas classes (normal e ataque), causou impacto direto na métrica de precisão, que se manteve entre 58,83% e 99,96%. Vale notar que no cálculo da precisão são considerados apenas os fluxos classificados como ataque. Diferentemente, a métrica de acurácia considera as duas classes (normal e ataque) em seu cálculo, justificando assim a discrepância entre os valores dessas duas medidas de desempenho. A melhoria nos resultados que os atributos avançados (entropias e demais medidas estatísticas) proporcionaram ao método *Bagging* pode ser atribuída ao fato do método em questão usar mais modelos de classificação do que o método VHT, que usa apenas um. Isso faz com que o mesmo consiga capturar com maior exatidão a contribuição, em termos de ganho de informações, de cada conjunto de atributos assinalados aos seus modelos de classificação. Assim sendo, no passo de decisão da classificação final de um fluxo, onde existe uma votação entre os modelos, fará diferença escolher aqueles atributos com maior contribuição para a formação do modelo de classificação.

Já os resultados ótimos a partir da utilização do método *Adaptive Bagging*, para o Grupo 3, podem ser explicados em função do mecanismo de funcionamento do método. Trata-se de um método mais otimizado para fluxos que mudam suas características ou padrões comportamentais ao longo do tempo. Na fase de construção do modelo, a acurácia é monitorada e, caso exista alguma pequena queda no valor da mesma, o modelo é refeito utilizando como base os últimos fluxos analisados. Dessa forma, tende-se a construção de um modelo bem ajustado aos padrões comportamentais dos fluxos, mesmo quando se usa um conjunto menor de atributos (no caso, os atributos primários). Entretanto, vale ressaltar que, do ponto de vista genérico, a estratégia usada por este método nem sempre pode ser tão eficiente, pois depende da dinâmica de variações entre as classes ao longo do *dataset* de treinamento.

Objetivando fornecer uma leitura mais fluida dos resultados alcançados, os gráficos da Figura 2 resumiram-nos para as métricas Precisão, *Recall* e *F-Measure* aplicadas ao Grupo 3. Nela, cada coluna representa o mecanismo de aprendizado de máquina utilizado (VHT, *Bagging* ou *Adaptive Bagging*) nos cenários CEN 4, CEN 10 e CEN 11 de modo a exibir o impacto nas métricas de desempenho apresentadas em função da variação dos conjuntos de atributos (Primários, Primários+entropias, Primários+estatísticos ou Primários+entropias+estatísticos).

Quanto às considerações de ordem operacional, o sistema consumiu, em média, 15 segundos com o método VHT, 17 segundos com o *Bagging* e 19 segundos com o *Adaptive Bagging*. Obteve-se um tempo médio de 42 microssegundos para a classificação de mais de 407.242 fluxos com 34 atributos cada um. Tais valores para o tempo de execução sugerem que, do ponto de vista da operação do sistema proposto em regime de produção num ambiente real, é possível a realização de uma etapa de treinamento extremamente rápida e capaz de gerar um sistema de detecção altamente acurado e com baixo tempo de reação. Vale ressaltar que os tempos médios calculados

já incluem o tempo gasto com a construção do modelo de classificação.

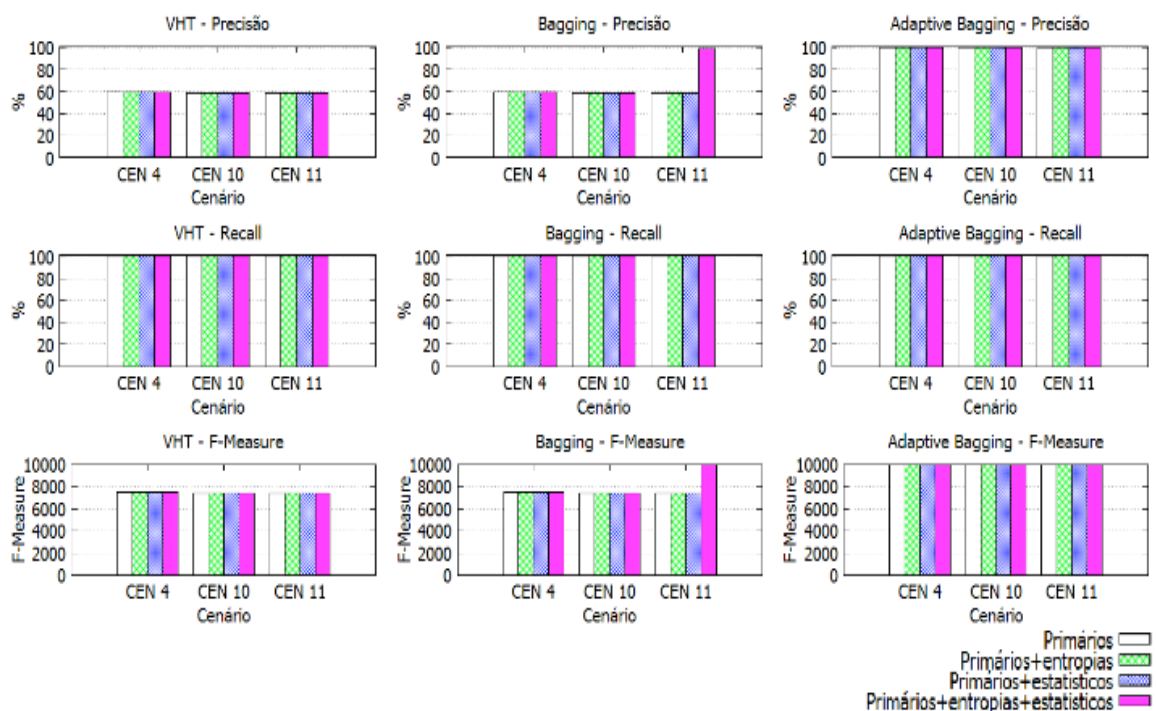


Figura 2. Resumo dos resultados para o Grupo 3

7 | CONCLUSÃO

Detectar com eficiência essas atividades maliciosas numa rede de computadores é uma tarefa cada vez mais complexa que se faz necessária pela importância das mesmas nas nossas vidas. Nesse contexto, este artigo analisa o impacto que determinados atributos, associados ao conjunto de fluxos de dados que trafegam por uma rede, pode exercer sobre mecanismos de aprendizado de máquina usados para a detecção de ataques. Para tal, foi desenvolvido um sistema de detecção de ataques DDoS baseado na arquitetura *lambda* e contendo mecanismos de aprendizado de máquina. A alimentação desse sistema foi efetuada via *datasets* realísticos gerados de fontes distintas, de forma a dificultar o processo de descoberta de conhecimento dos padrões comportamentais dos fluxos de dados que passam pela rede. Os resultados mostram que a variação do conjunto de atributos associado a esses fluxos pode causar impacto nas métricas de desempenho, especialmente em casos onde há variação considerável na natureza desses fluxos. Verificou-se que é possível aprimorar a classificação dos fluxos de dados realizada com o método *Bagging* quando se incorpora, ao conjunto de atributos, medidas estatísticas que refletem a dinâmica comportamental do tráfego agregado da rede. Além disso, o sistema desenvolvido se mostrou suficientemente rápido e preciso para uso em cenários reais de operação.

Como trabalhos futuros, deseja-se ampliar a avaliação do sistema para a detecção de outros tipos de ataques que explorem a camada de aplicação. Deseja-se, também, identificar novos atributos estatísticos relacionados à dinâmica do tráfego

agregado, que porventura melhorem o desempenho dos mecanismos de aprendizado de máquina na identificação dos ataques. Por fim, deseja-se implementar o sistema de detecção proposto num ambiente de produção.

REFERÊNCIAS

BHUYAN, Monowar H.; BHATTACHARYYA, Dhruva Kumar; KALITA, Jugal K. Network anomaly detection: methods, systems and tools. **IEEE communications surveys & tutorials**, v. 16, n. 1, p. 303-336, 2014.

DE FRANCISCI MORALES, Gianmarco. SAMOA: a platform for mining big data streams. In: **Proceedings of the 22nd International Conference on World Wide Web**. ACM, 2013. p. 777-778.

DU, Yutan et al. A real-time anomalies detection system based on streaming technology. In: **2014 Sixth International Conference on Intelligent Human-Machine Systems and Cybernetics**. IEEE, 2014. p. 275-279.

Stealthy denial of service strategy in cloud computing

GARCIA, Sebastian et al. An empirical comparison of botnet detection methods. **computers & security**, v. 45, p. 100-123, 2014.

GARCIA, S.; UHLIR, V. (2017). Malware capture facility project - ctu-13 data-set. <http://mcfp.weebly.com/te-ctu-13-dataset-a-labeled-dataset-with-botnet-normal-and-background-traffic.html>. Acessado em 07/04/2017.

JOLDZIC, Ognjen; DJURIC, Zoran; VULETIC, Pavle. A transparent and scalable anomaly-based DoS detection method. **Computer Networks**, v. 104, p. 27-42, 2016.

KHAN, Latifur; AWAD, Mamoun; THURASINGHAM, Bhavani. A new intrusion detection system using support vector machines and hierarchical clustering. **The VLDB journal**, v. 16, n. 4, p. 507-521, 2007.

LOBATO, A.; ANDREONI LOPEZ, M.; DUARTE, O. C. M. B. Um sistema acurado de detecção de ameaças em tempo real por processamento de fluxos. **XXXIV Simpósio Brasileiro de Redes de Computadores e Sistemas Distribuídos (SBRC'2016)**, Salvador, Bahia, 2016.

MANIKANDAN, Shankar Ganesh; RAVI, Siddarth. Big data analysis using Apache Hadoop. In: **2014 International Conference on IT Convergence and Security (ICITCS)**. IEEE, 2014. p. 1-4.

MARZ, Nathan; WARREN, James. **Big Data: Principles and best practices of scalable real-time data systems**. New York; Manning Publications Co., 2015.

MICHIE, Donald et al. Machine learning. **Neural and Statistical Classification**, v. 13, 1994.

QUINLAN, J.. Ross . Induction of decision trees. **Machine learning**, v. 1, n. 1, p. 81-106, 1986.

ROBINSON, RR Rejimol; THOMAS, Ciza. Ranking of machine learning algorithms based on the performance in classifying ddos attacks. In: **2015 IEEE Recent Advances in Intelligent Computational Systems (RAICS)**. IEEE, 2015. p. 185-190.

SHANNON, Claude Elwood. A mathematical theory of communication. **Bell system technical journal**, v. 27, n. 3, p. 379-423, 1948.

SINGH, Kamaldeep et al. Big data analytics framework for peer-to-peer botnet detection using random forests. **Information Sciences**, v. 278, p. 488-497, 2014.

TOSHNIWAL, Ankit et al. Storm@ twitter. In: **Proceedings of the 2014 ACM SIGMOD international conference on Management of data**. ACM, 2014. p. 147-156.

YI, Xiaomeng et al. Building a network highway for big data: architecture and challenges. **IEEE Network**, v. 28, n. 4, p. 5-13, 2014.

SOBRE O ORGANIZADOR

Ernane Rosa Martins - Doutorado em andamento em Ciência da Informação com ênfase em Sistemas, Tecnologias e Gestão da Informação, na Universidade Fernando Pessoa, em Porto/Portugal. Mestre em Engenharia de Produção e Sistemas, possui Pós-Graduação em Tecnologia em Gestão da Informação, Graduação em Ciência da Computação e Graduação em Sistemas de Informação. Professor de Informática no Instituto Federal de Educação, Ciência e Tecnologia de Goiás - IFG (Câmpus Luziânia), ministrando disciplinas nas áreas de Engenharia de Software, Desenvolvimento de Sistemas, Linguagens de Programação, Banco de Dados e Gestão em Tecnologia da Informação. Pesquisador do Núcleo de Inovação, Tecnologia e Educação (NITE), certificado pelo IFG no CNPq.

Agência Brasileira do ISBN
ISBN 978-85-7247-478-8

