

Classificação de Tráfego Darknet

Universidade Federal do Paraná
CI1030 Ciência de Dados para Segurança

Jackson Rossi Borguezani



Objetivo

Detectar o tráfego darknet, combinando a categorização de dois conjuntos de dados públicos, o tráfego Tor, Non-Tor e VPN, Non-VPN.

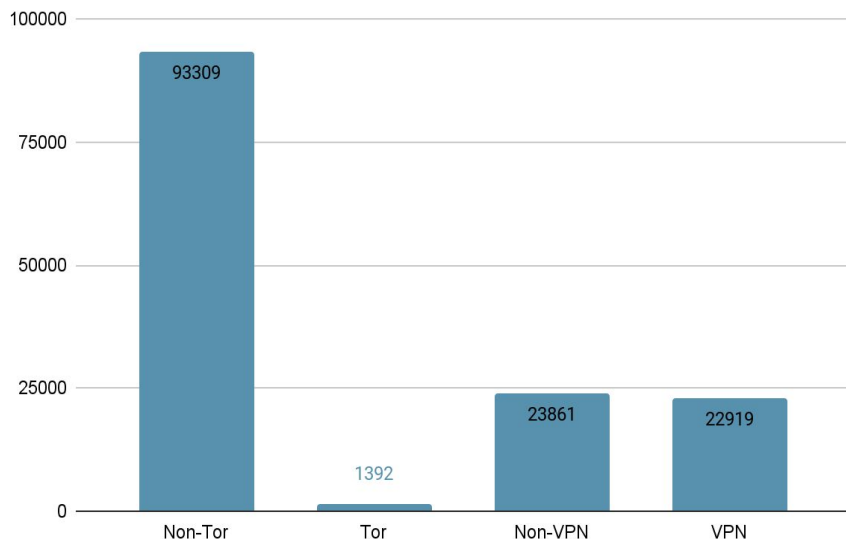


Introdução

- CIC-Darknet2020 [1]
- Arquivo usado: Darknet.csv
- 85 colunas combinadas, contendo dados do tráfego VPN e Tor



Exploração de Dados



- Dados majoritariamente numéricos
- Má distribuição dos dados de tráfego Tor, Non-Tor
- 31 colunas pouco relevantes foram descartadas nesta análise



Extração de Características

- Matriz de correlação entre 25 colunas
- Remoção de 8 colunas correlacionados que podem ser ignorados



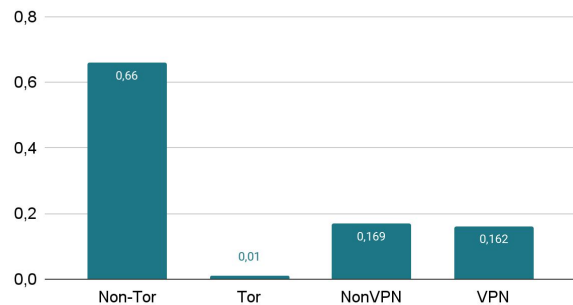
Machine Learning

1. Separar o dataset em 80% e 20%
2. Treinar com a parte de 80%
 - a. Validar modelos com split 80/20
 - b. Validar modelos com KFold de 5
3. Testar com os 20%

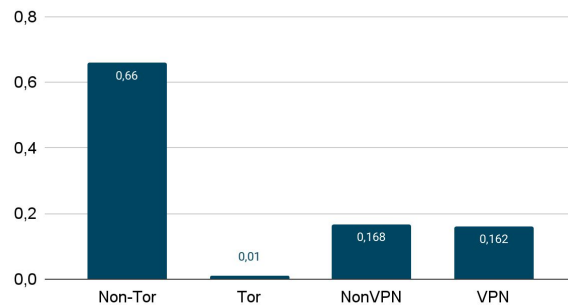


1. Separação do dataset

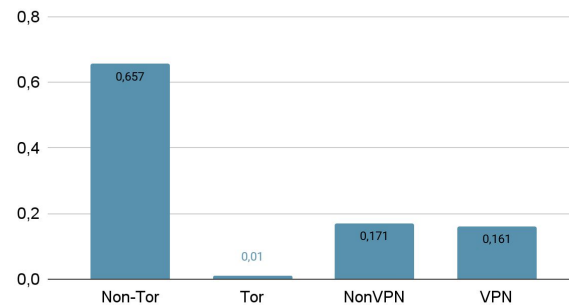
Dataset Inteiro



Dataset 80%



Dataset 20%



2. Treinamento

- Normalização usando MinMaxScaler
- Foram treinados 4 modelos

SVM

class_weight	'balanced'
random_state	21
cache_size	500
probability	True
max_iter	1000

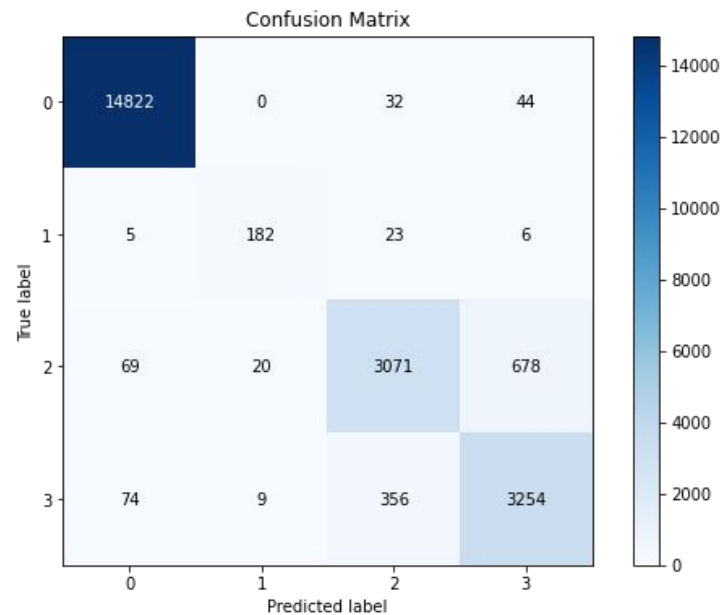
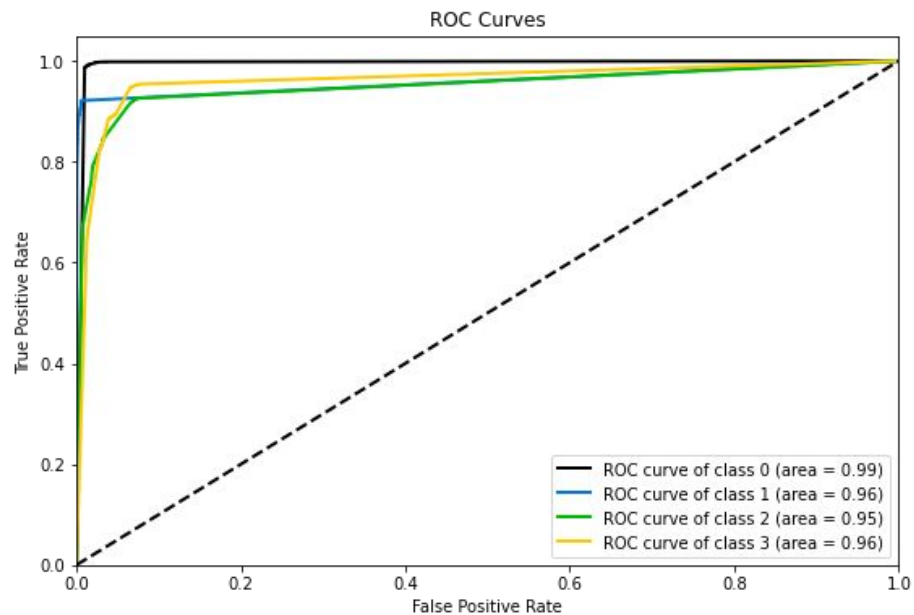
KNN

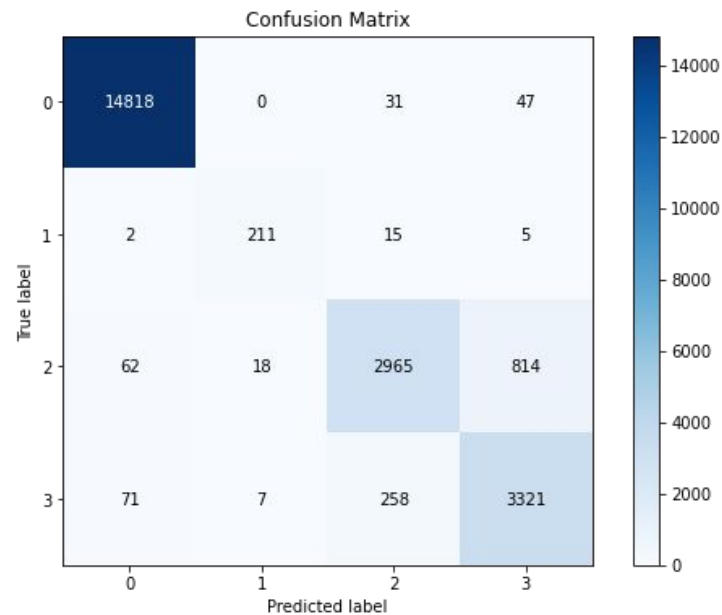
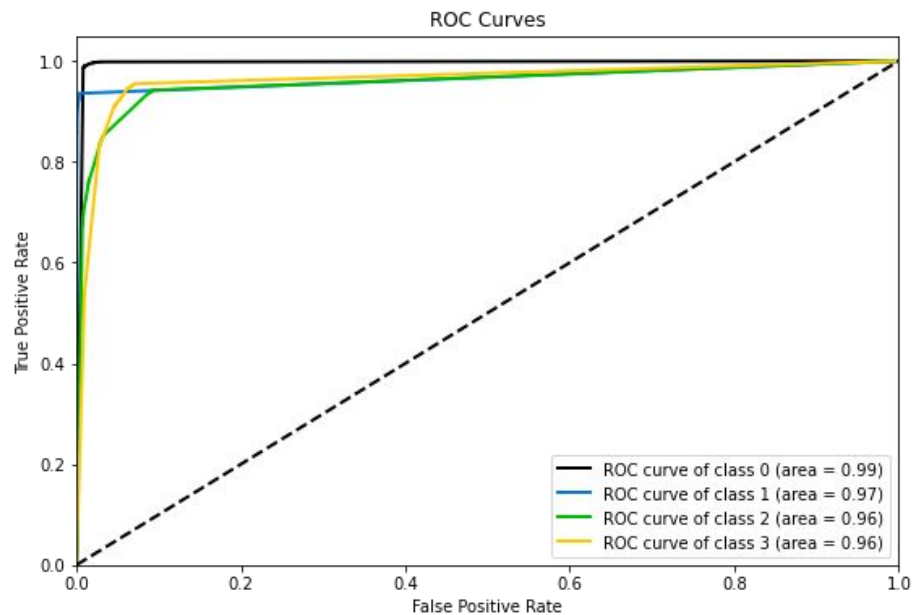
n_neighbors	4
weights	'distance'

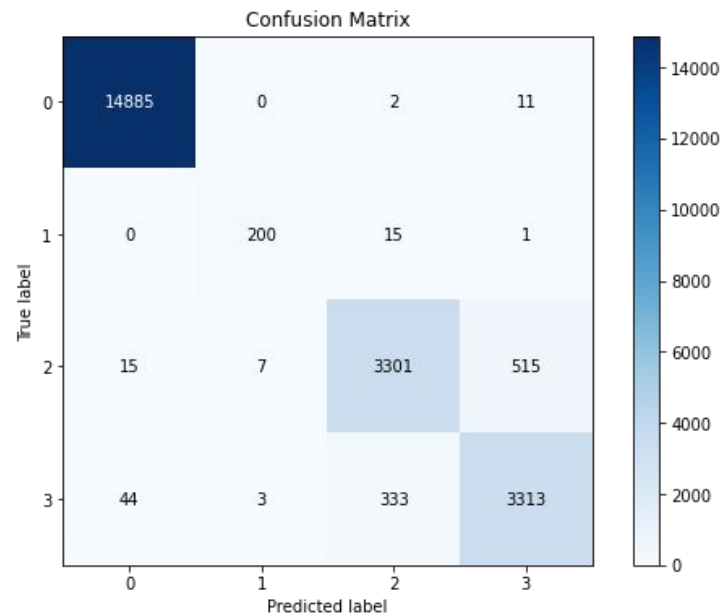
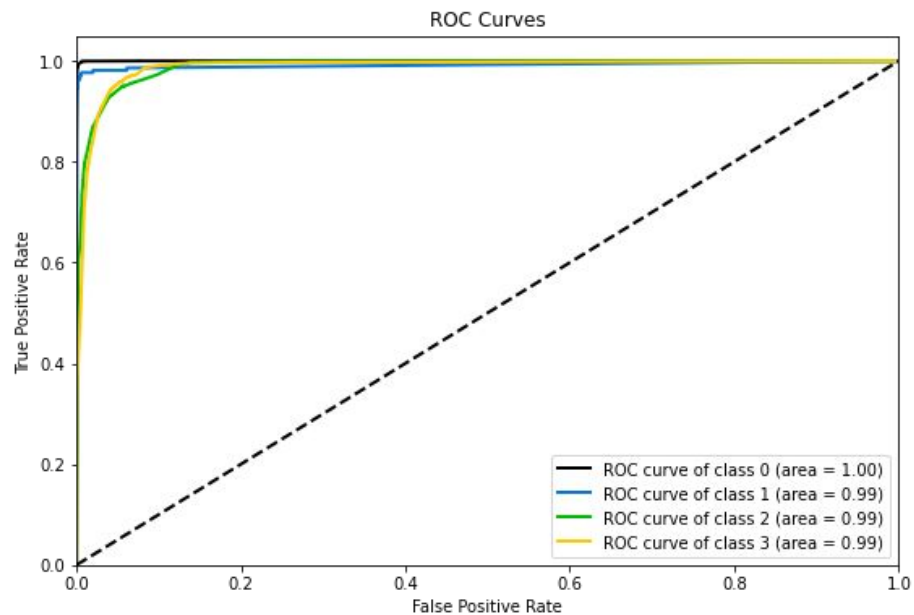
Random Forest

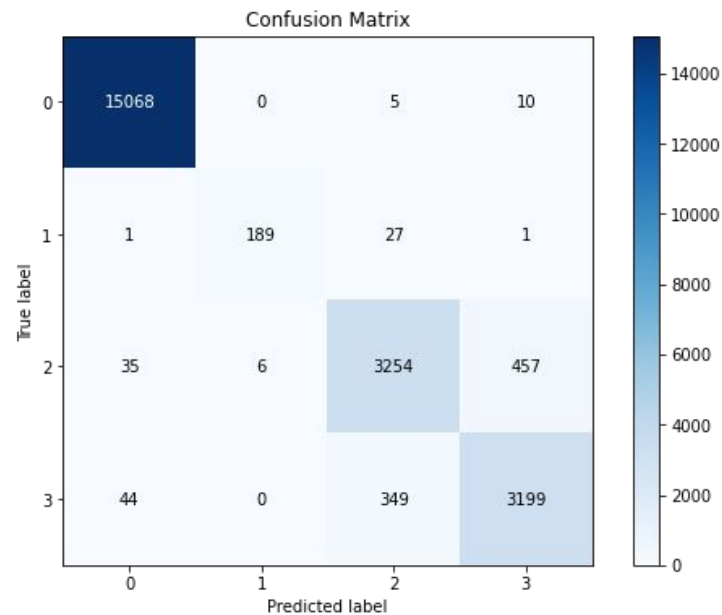
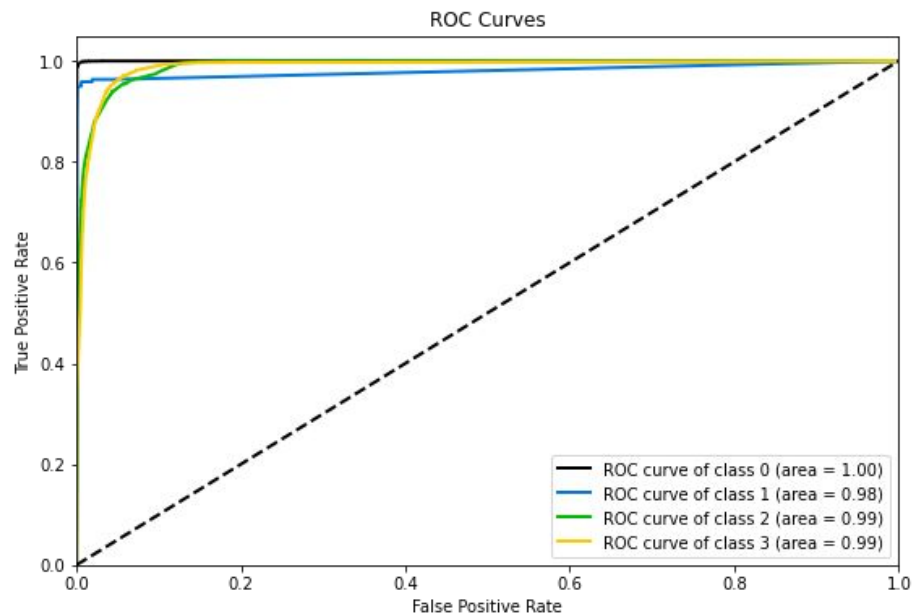
n_estimators	75
random_state	21

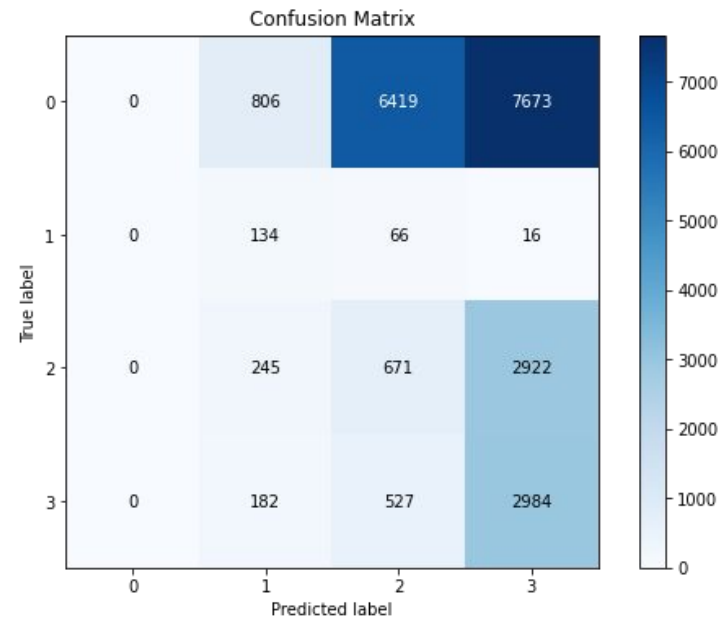
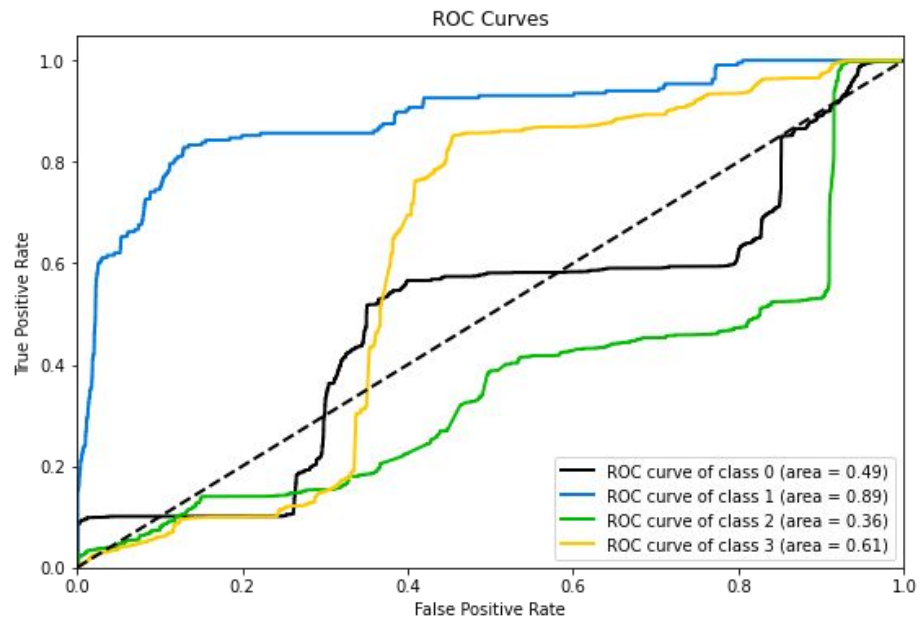


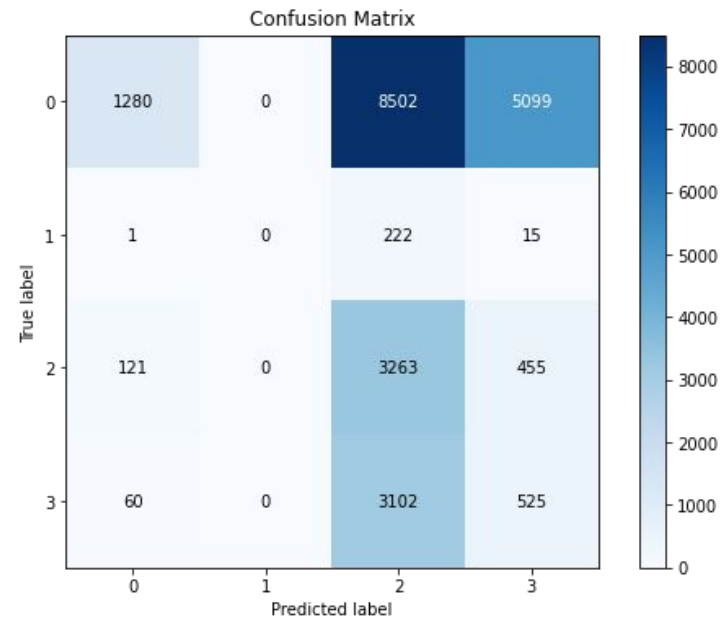
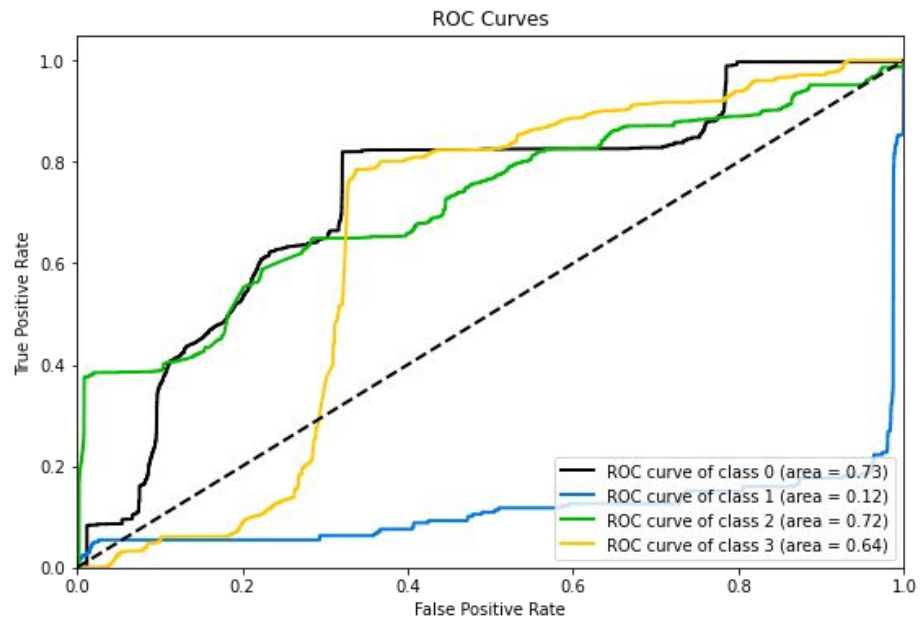






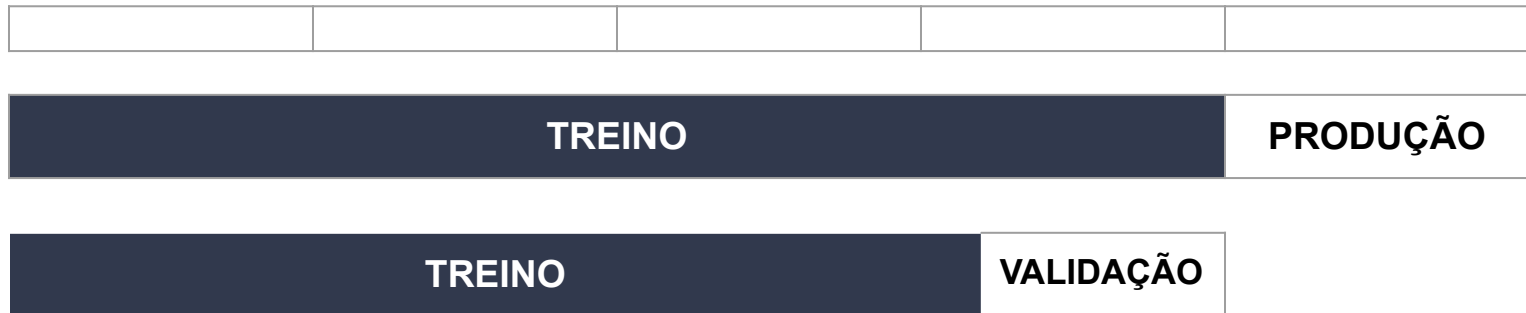


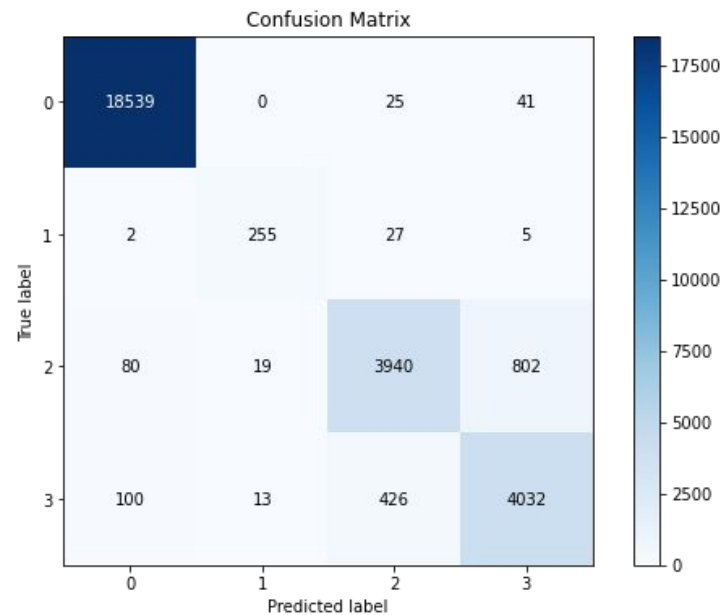
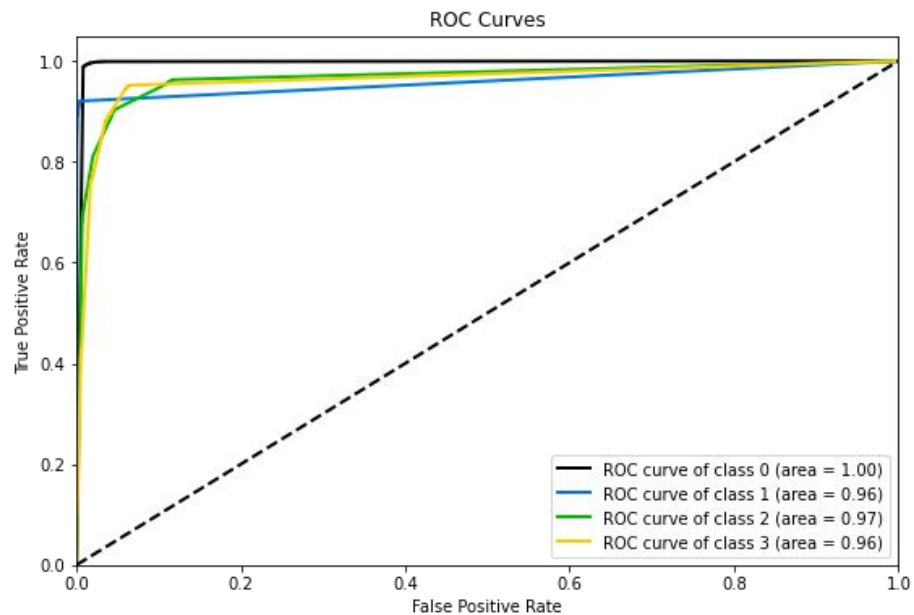


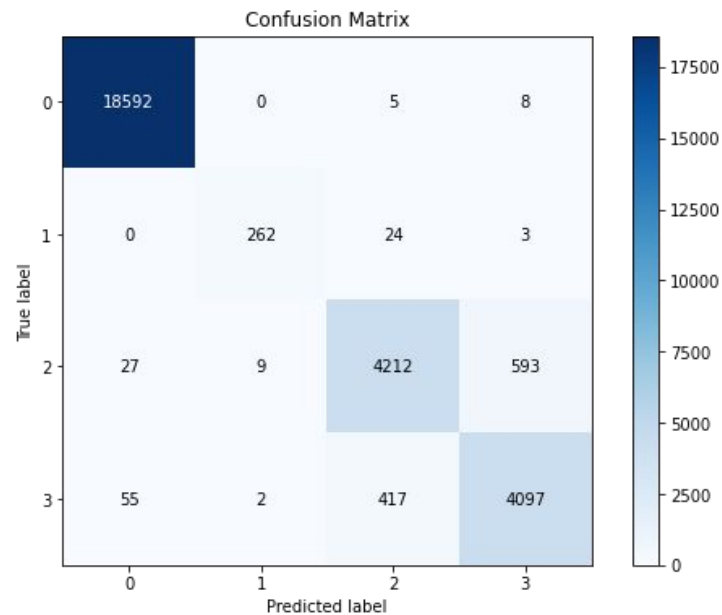
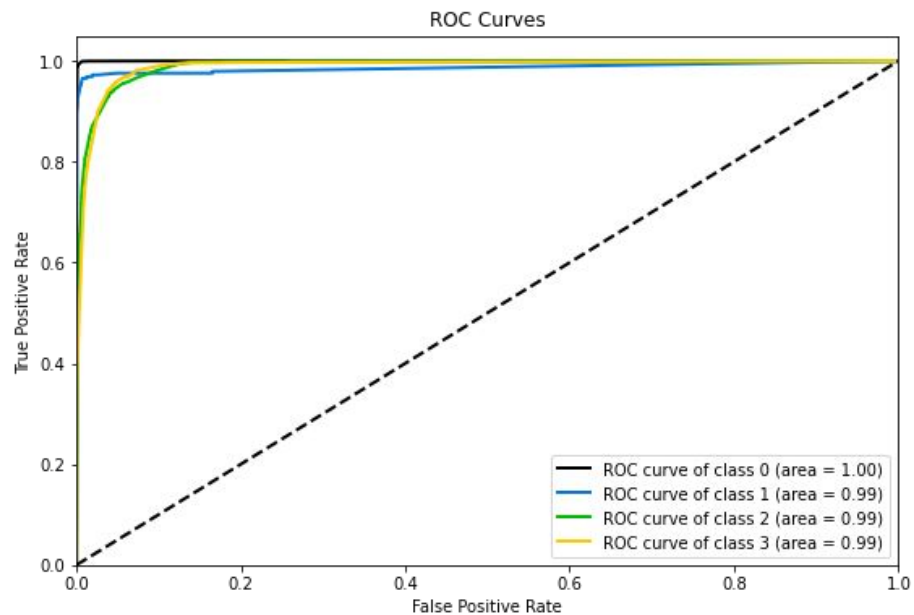


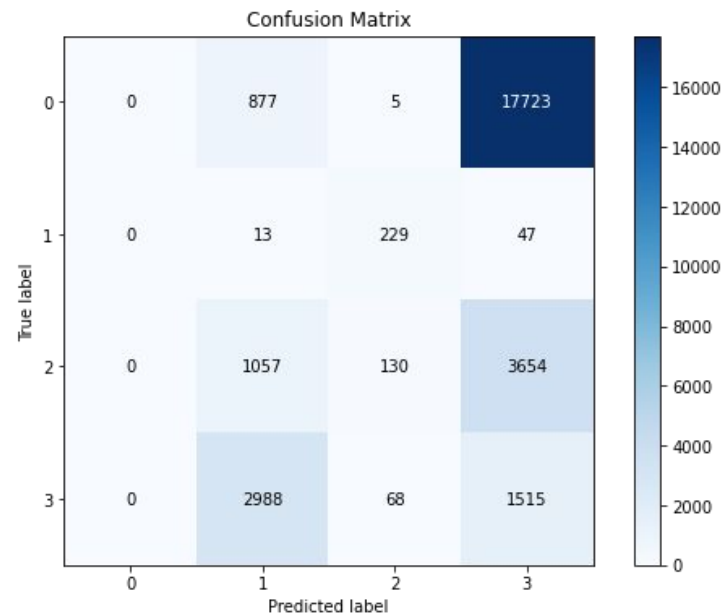
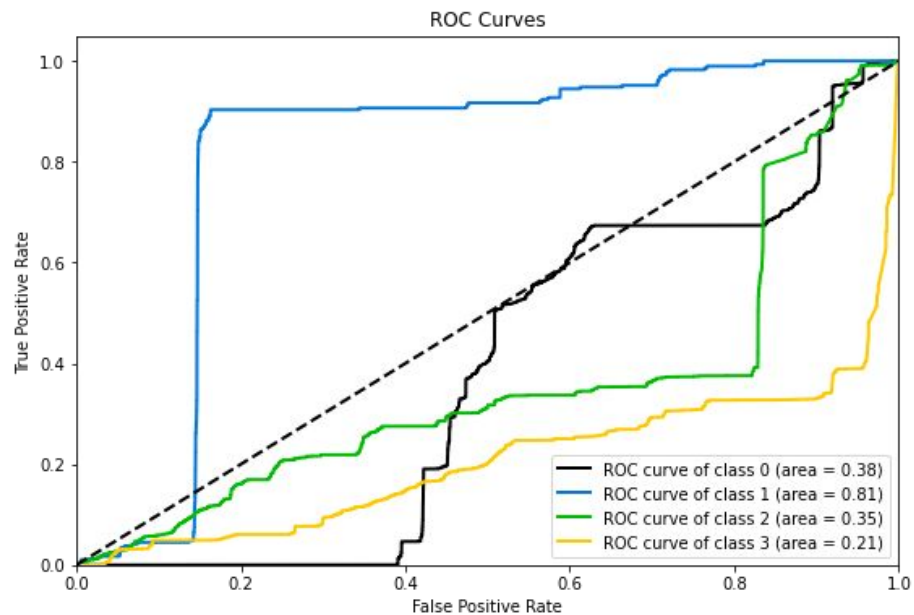
3. Teste

- Usando a parte de 20% como dados de produção
- Aplicado aos modelos treinados com a porção inteira de 80%

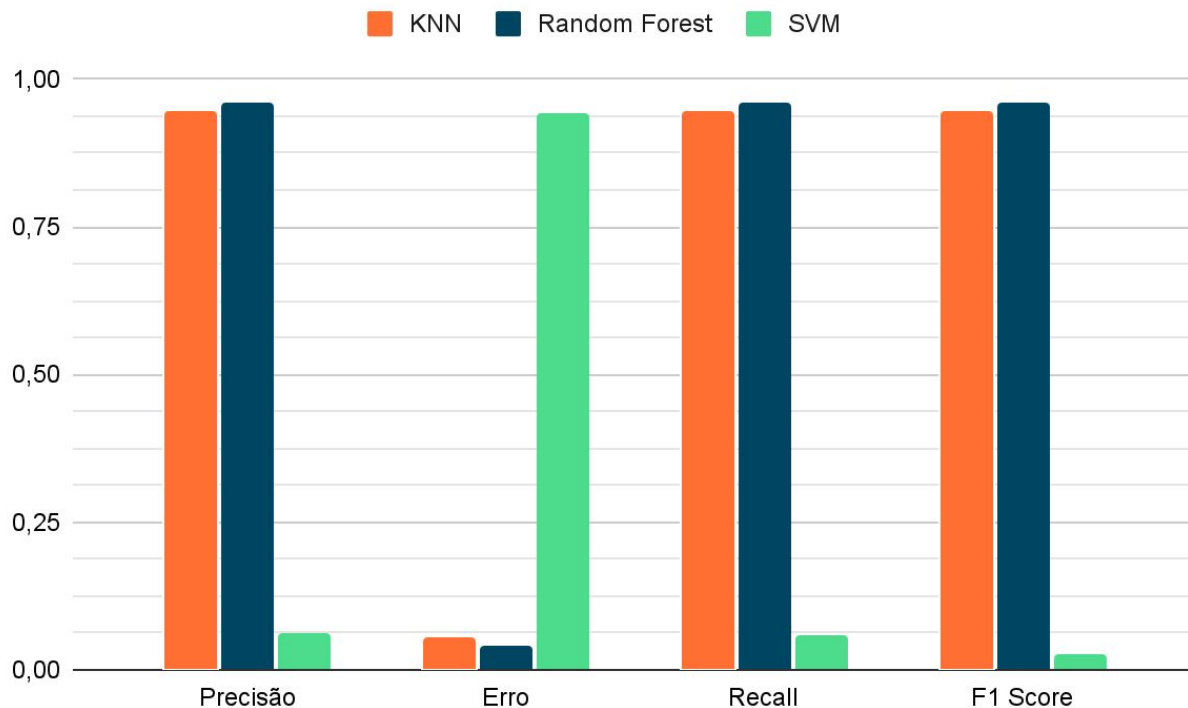








Resultado dos Testes



Conclusão

De modo geral, o Random Forest parece a melhor solução para classificar o tráfego:

- Tem a melhor precisão
- Menor erro
- Menor tempo de predição

O SVM teve um desempenho baixíssimo, já que seu tempo de treino foi muito grande e precisou de limitação de iterações. O fato do dataset ser multiclasse pode ser um agravante no tempo de treino, além do tamanho do dataset. Então o SVM não é viável de forma alguma para esses dados.

Já o KNN teve um tempo de treino melhor que o Random Forest. No entanto, a penalidade acontece no tempo de predição, que é muito maior ao Random Forest.

Assim, o Random Forest é o melhor algoritmo de machine learning entre os três estudados.



Links e Referências

- GitHub: <https://github.com/jacksonrossi/ciencia-dados-darknet>
- [1] Arash Habibi Lashkari, Gurdip Kaur, and Abir Rahali, “DIDarknet: A Contemporary Approach to Detect and Characterize the Darknet Traffic using Deep Image Learning”, 10th International Conference on Communication and Network Security, Tokyo, Japan, November 2020. <https://www.unb.ca/cic/datasets/darknet2020.html>



