

CI1030 – Ciência de Dados para Segurança

Exploração e classificação de um dataset de URLs maliciosas

Aluno: Christian Debovi Paim Oliveira (GRR20186713)
Professor: André Gregio

Dataset – Pesquisa

- **Paper:** Detecting Malicious URLs Using Lexical Analysis.
- **Autores:**
 - Mohammad Saiful Islam Mamun
 - Mohammad Ahmad Rathore
 - Arash Habibi Lashkari
 - Natalia Stakhanova
 - Ali A. Ghorbani
- **Universidade:**
 - University of New Brunswick , Fredericton, NB, Canada
- **Objetivo:**
 - Avaliar o uso de características léxicas para a classificação de URLs
 - Avaliar uso de técnicas de ofuscação nas urls coletadas

Datasets – URLs

Arquivos texto com urls:

- Total de 165.366 URLs

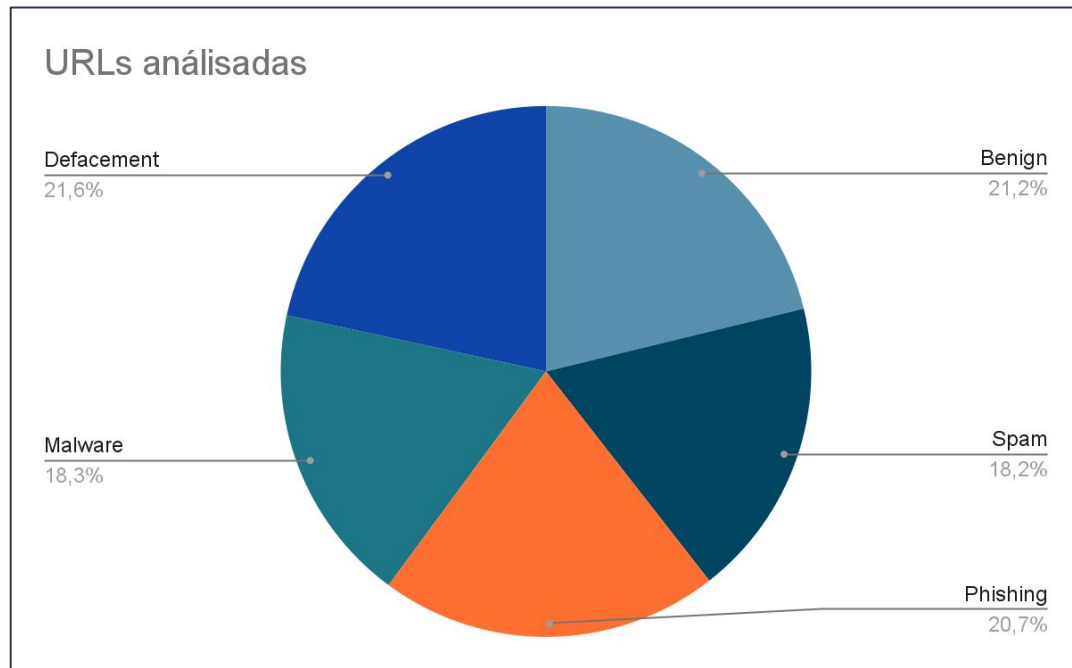
Labels/Classes:

- **Benign:** Alexa top sites + crawler + VirusTotal (35.378 URLs).
- **Spam:** dataset WEBSPPAM-UK2007 (12.000 URLs).
- **Phishing:** repositório OpenPhish (9965 URLs).
- **Malware:** lista DNS-BH (11.566 URLs).
- **Defacement:** em Alexa top sites (96.457 URLs).

Dataset – Informações de URLs

CSVs com informações léxicas:

- Extraídas de 36.707 URLs do total
- 79 atributos léxicos + label



Dataset – Atributos

Tipos:

- **Entropy:** variação nos tokens em certas partes da url.
 - Entropy_Domain, Entropy_Extension
- **CharacterContinuityRate:** soma da continuidade dos tokens divididos pelo tamanho da URL.
 - **Ex:** $abc567ti = (3 + 3 + 1)/9 = 0.77$
- **Ratios:** número de tokens de uma parte da URL divididos pelo número de tokens da outra.
 - argPathRatio, argUrlRatio, argDomainRatio, domainUrlRatio, pathUrlRatio, PathDomainRatio.
- **NumberRate:** proporção de dígitos nas partes da URL.
 - NumberRate_Domain, NumberRate_DirectoryName, NumberRate_FileName, NumberRate_URL, NumberRate_AfterPath.
- Outros atributos relacionados ao tamanho e contagem de tokens, dígitos e símbolos de diferentes partes da URL.

Datasets – Seleção de Atributos (WEKA)

Infogain:

Determina o “peso” de um atributo através da medição do ganho de informação em respeito à classe

$\text{InfoGain}(\text{Class}, \text{Attribute}) =$

$H(\text{Class}) - H(\text{Class} \mid \text{Attribute}).$

Ranker

Ordena atributos de acordo com sua avaliação.

The screenshot shows the Weka Explorer window with the 'Select attributes' tab selected. The 'Attribute Evaluator' is set to 'InfoGainAttributeEval' and the 'Search Method' is set to 'Ranker'. The 'Attribute selection output' pane displays the results of the attribute selection process, including the ranked attributes and the search method used.

Attribute selection output

=== Attribute Selection on all input data ===

Search Method:
Attribute ranking.

Attribute Evaluator (supervised, Class (nominal): 80 URL_Type_obf_Type):
Information Gain Ranking Filter

Ranked attributes:

1.590258	75	Entropy_Domain
1.038767	32	argPathRatio
1.023913	28	argUnRatio
1.018038	29	argDomainRatio
0.964314	27	pathUnRatio
0.949599	37	CharacterContinuityRate
0.906596	65	NumberRate_FileName
0.896363	30	domainUnRatio
0.876959	62	NumberRate_URL
0.875926	31	pathDomainRatio
0.721321	67	NumberRate_AfterPath
0.706483	6	avgpathTokenlen
0.673254	78	Entropy_Extension
0.653202	4	avgdomainTokenlen
0.642828	79	Entropy_Afterpath
0.602251	77	Entropy_FileName
0.575901	51	LongestPathTokenLength
0.572804	76	Entropy_DirectoryName
0.548466	38	LongestVariableValue
0.520601	35	NumbersDotsinURL
0.51097	23	subDirLen
0.51097	22	pathLength
0.510708	20	urlLen
0.507783	26	argLen
0.498566	21	domainlength
0.475868	66	NumberRate_Extension
0.474437	1	Querylength
0.471897	50	Query_LetterCount
0.463078	49	Extension_LetterCount

Dataset – Atributos escolhidos

Escolhidos 12 de 79 atributos

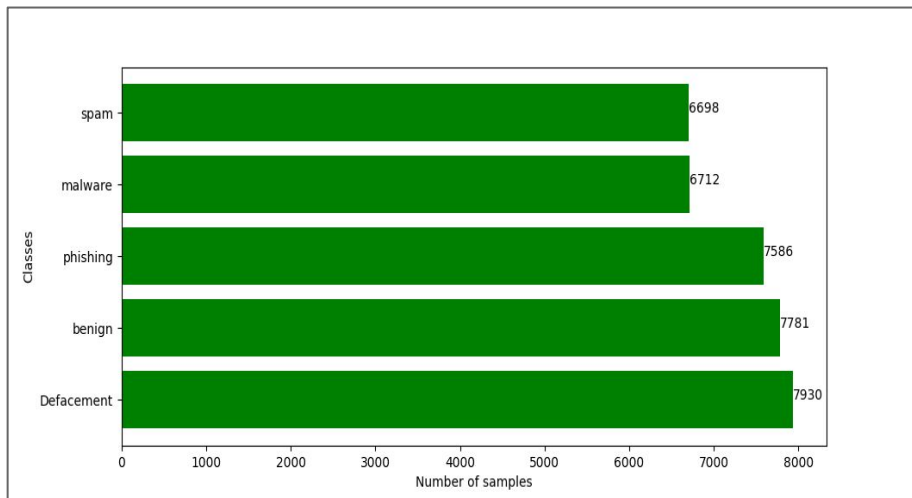
- Entropy_Domain
- argPathRatio
- ArgUrlRatio
- argDomanRatio
- pathurlRatio
- CharacterContinuityRate
- NumberRate_FileName
- domainUrlRatio
- NumberRate_URL
- pathDomainRatio
- NumberRate_AfterPath
- avgpathtokenlen

Dataset – Atributos escolhidos

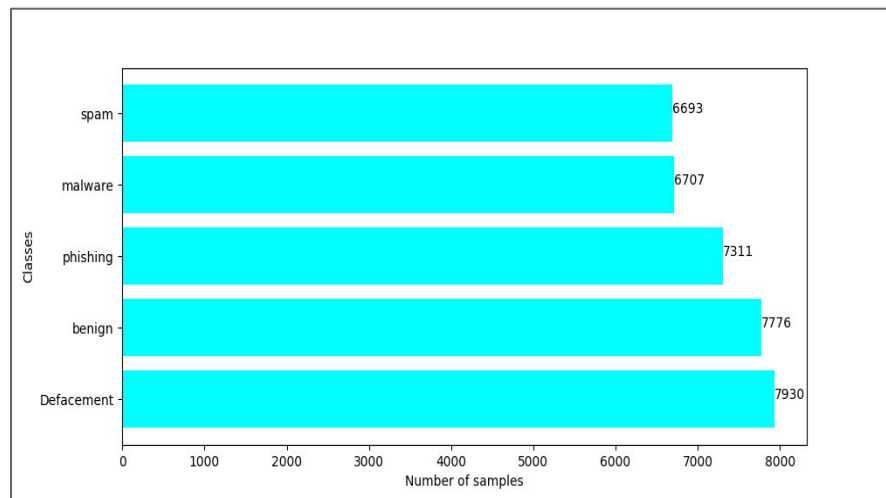
[illegible]

Dataset – Distribuição

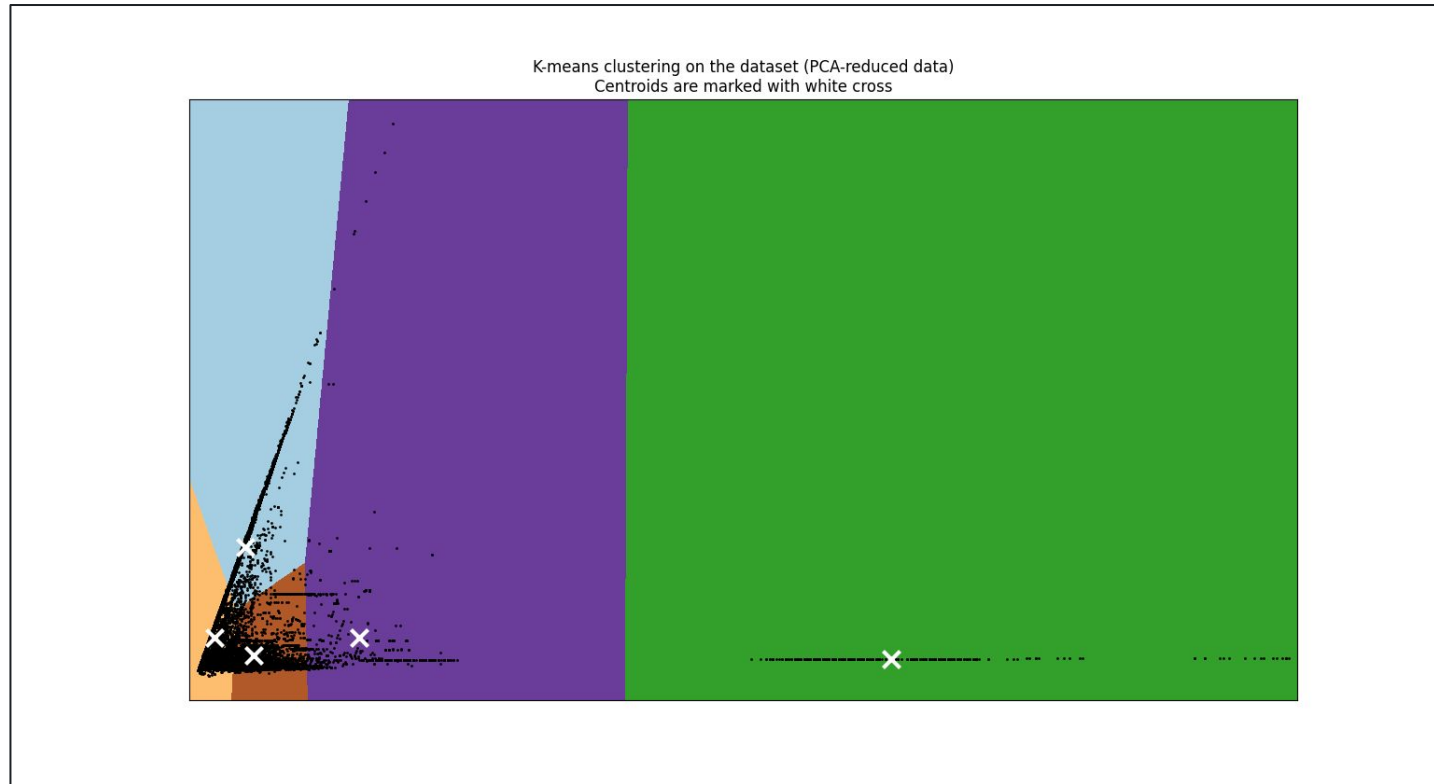
Todos os dados



Infogain + Dropnan



Dataset – Clusters



Fontes

- Dataset utilizado:
 - <https://www.unb.ca/cic/datasets/url-2016.html>
- Paper “Detecting Malicious URLs Using Lexical Analysis”:
 - https://www.researchgate.net/publication/308365207_Detecting_Malicious_URLs_Using_Lexical_Analysis

Obrigado pela atenção!

