

# Investigando o uso de Características na Detecção de URLs Maliciosas

**Gustavo Ferreira da Silva**

# URL

- Sigla para Uniform Resource Locator.
- Composta da seguinte forma:
- <esquema>:<parte-especifica-do-esquema>

# Composição da URL

• `http://`    `icomp.ufam`    `.edu.`    `br`



Esquema



Nome do Domínio  
Secondary Level  
Domain (SLD)



Generic  
Top  
Level  
Domain  
(GTLD)

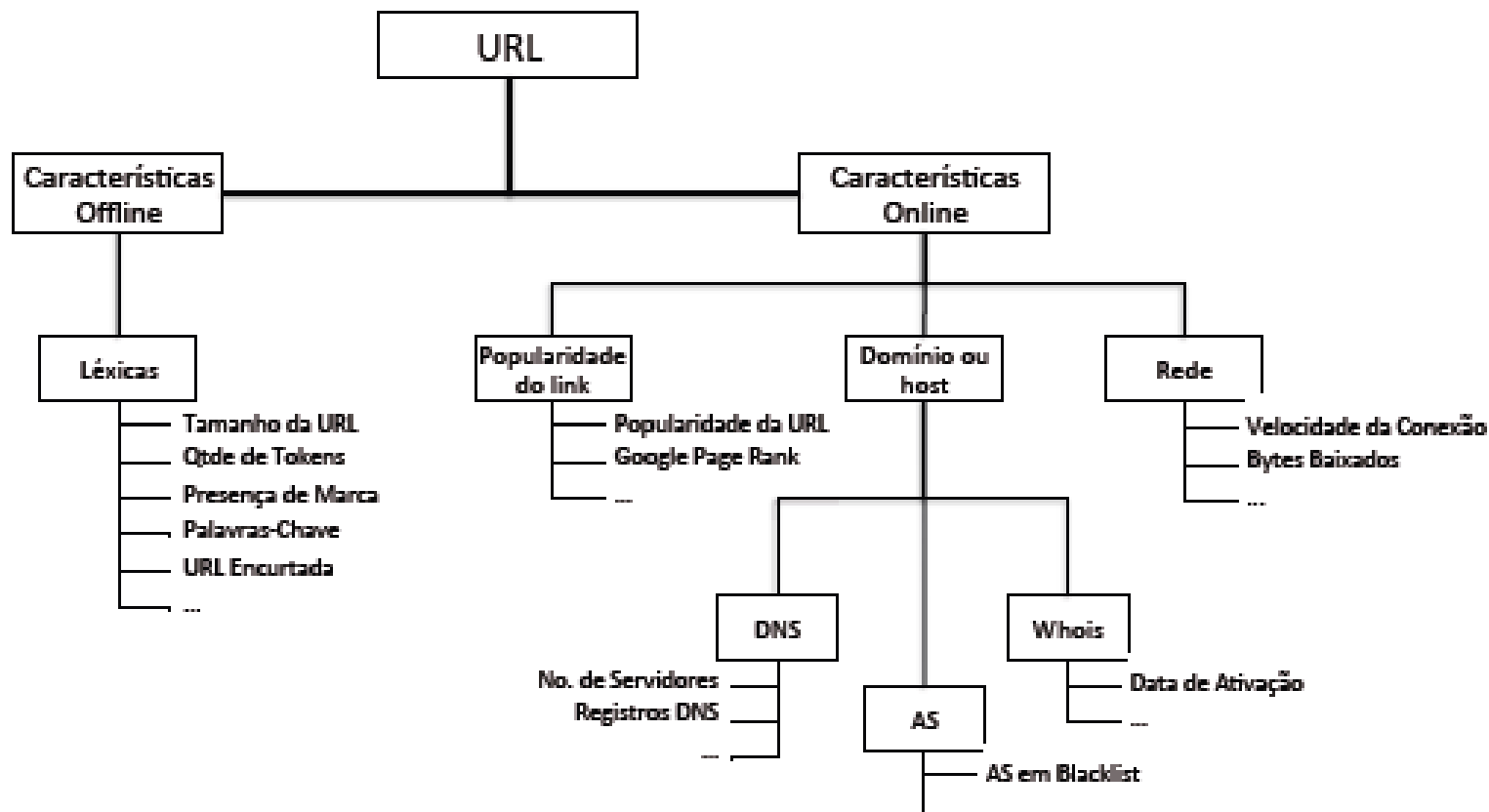


Country  
Code  
(CCTLD)

# Componentes de uma URL

Componente	Exemplo
URL	http://icomp.ufam.edu.br/inst/hstart.php?id=664&logon=141
Nome do domínio	icomp.ufam.edu.br
Caminho	inst/hstart.php
Sub diretório	inst
Nome do arquivo	hstart
Extensão do arquivo	php
Argumento	id=664&logon=141

# Agrupamento das URLs



# Algoritmos de Aprendizagem de Máquina

- SVM
- Naive Bayes
- Árvore de Decisão (J.48)
- KNN

# Características Léxicas

Características Léxicas					
Nome	Descrição	Nome	Descrição	Nome	Descrição
<i>qt_dom_ponto</i>	Qtde de (.) no domínio	<i>qt_dom_hifen</i>	Qtde de (-) no domínio	<i>qt_dom_underline</i>	Qtde de ( _ ) no domínio
<i>qt_url_ponto</i>	Qtde de (.) na URL	<i>qt_url_barra</i>	Qtde de (/) na URL	<i>qt_url_interrog</i>	Qtde de (?) na URL
<i>qt_url_igualdade</i>	Qtde de (=) na URL	<i>qt_url_hifen</i>	Qtde de (-) na URL	<i>qt_url_underline</i>	Qtde de ( _ ) na URL
<i>qt_url_arroba</i>	Qtde de (@) na URL	<i>qt_url_comerc</i>	Qtde de (&) na URL	<i>qt_url_exclam</i>	Qtde de (!) na URL
<i>qt_url_til</i>	Qtde de () na URL	<i>comp dominio</i>	Comprimento do domínio	<i>comp_url</i>	Comprimento da URL
<i>qt_dir_ponto</i>	Qtde de (.) no diretório	<i>qt_dir_barra</i>	Qtde de (/) no diretório	<i>qt_dir_interrog</i>	Qtde de (?) no diretório
<i>qt_dir_igualdade</i>	Qtde de (=) no diretório	<i>qt_dir_hifen</i>	Qtde de (-) no diretório	<i>qt_dir_underline</i>	Qtde de ( _ ) no diretório
<i>qt_dir_arroba</i>	Qtde de (@) no diretório	<i>qt_dir_exclam</i>	Qtde de (!) no diretório	<i>qt_dir_til</i>	Qtde de () no diretório
<i>qt_arq_ponto</i>	Qtde de (.) no arquivo	<i>qt_arq_interrog</i>	Qtde de (?) no arquivo	<i>qt_arq_igualdade</i>	Qtde de (=) no arquivo
<i>qt_arq_hifen</i>	Qtde de (-) no arquivo	<i>qt_arq_underline</i>	Qtde de ( _ ) no arquivo	<i>qt_arq_arroba</i>	Qtde de (@) no arquivo
<i>qt_arq_exclam</i>	Qtde de (!) no arquivo	<i>qt_arq_til</i>	Qtde de () no arquivo	<i>qt_par_ponto</i>	Qtde de (.) no parâmetro
<i>qt_par_barra</i>	Qtde de (/) no parâmetro	<i>qt_par_interrog</i>	Qtde de (?) no parâmetro	<i>qt_par_igualdade</i>	Qtde de (=) no parâmetro
<i>qt_par_hifen</i>	Qtde de (-) no parâmetro	<i>qt_par_underline</i>	Qtde de ( _ ) no parâmetro	<i>qt_par_arroba</i>	Qtde de (@) no parâmetro
<i>qt_par_comerc</i>	Qtde de (&) no parâmetro	<i>qt_par_exclam</i>	Qtde de (!) no parâmetro	<i>qt_par_til</i>	Qtde de () no parâmetro
<i>qt_params</i>	Qtde de parâmetros na URL	<i>pres_tld_arg</i>	Presença de TLD no argumento da URL	<i>comp_dirtorio</i>	Comprimento do diretório da URL
<i>comp_arquivo</i>	Comprimento do arquivo na URL	<i>comp_params</i>	Comprimento dos parâmetros da URL	—	—

# Características de DNS

Nome	Descrição
<i>ip_associado</i>	No. de IPs resolvidos
<i>sn_associado</i>	No. de servidores de nome resolvidos
<i>data_tempo_ativo</i>	Tempo (em dias) de ativação do domínio



# Características Especiais

Nome	Descrição
<i>mal_phi</i>	Presença em listas de Phishing ou Malware
<i>rank_google</i>	Page Rank do Google
<i>presenca_marca</i>	Presença de marca
<i>rank_alexa</i>	Page Rank do Alexa
<i>geo_localizacao</i>	Localização geográfica do domínio
<i>rbl_check</i>	Presença do domínio em RBL ( <i>Real-time Blackhole List</i> )

# Ambiente Utilizado

- Windows 7 64 bits, 4 GB de memória RAM, disco de 500 GB e um processador Intel Core i5, 2.3 Ghz
- Intel Core 7 de 3.4 Ghz, com 8 GB de memória RAM, disco de 500 GB e plataforma Linux, distribuição Ubuntu 14.04
- Weka 3.6.10

# Bases de Dados Utilizadas

- DMOZ (URLs Benignas)
- PhishTank (URLs Maliciosas)
- Shalla's Blacklist (Agrupamentos de diversas URLs em categorias destinadas a filtros. Também destinada a URLs Maliciosas).

# Quantidade de URLs utilizadas

## Treinamento e Ajuste

- 20.092 – Total
- 10.046 – Base DMOZ
- 10.046 – Base PhishTank

# Quantidade de URLs utilizadas

## Teste

- 20.000 – Base DMOZ e PhishTank de forma igualitária
- 20.000 – Base DMOZ e Shella's forma igualitária

# Medidas de Desempenho

- Taxa de Detecção:  $VP / (VP+FN)$
- Taxa de Precisão:  $(VP+VN) / (VP+VN+FP+FN)$
- Taxa de Falso Alarme:  $FP / (FP+VN)$
  
- VN (Verdadeiro Negativo)
- VP (Verdadeiro Positivo)
- FP (Falso Positivo)
- FN (Falso Negativo)

# Métrica empregada

- 10 partições para os 4 classificadores.

# Comparação

Classificador	Naive Bayes	SVM	Árvore de Decisão	KNN
Parâmetros Ajustados	-	<i>C=50, Kernel Polinomial, Grau do Polinômio= 1.0</i>	Fator de Confiança= 0,25	KNN=1
Taxa de Precisão	76,00%	91,40%	95,10%	94,90%
Taxa de Detecção	66,35%	91,43%	95,11%	94,91%
Falso Alarme	33,60%	8,60%	4,90%	5,10%



# Hipóteses

- **H1. Existe alguma influência do formato da URL na extração das características e, conseqüentemente, no processo de avaliação?**
- **H2. Todas as características extraídas são realmente necessárias no processo de detecção de URLs?**
- **H3. Grupos de características permitem resultados adequados, e até melhores, no processo de detecção de URLs se comparados com características individuais.**
- **H4. A importância das características depende da base onde as URLs são coletadas.**

# Provas H1

Características	Shalla's (Dez. 2014)	PhishTank (Nov. 2013)
<i>comp_url</i>	31,0367	54,1183
<i>comp_dominio</i>	14,3585	19,0749
<i>comp_diretorio</i>	6,1998	21,6629
<i>comp_arquivo</i>	5,3814	7,1557
<i>comp_params</i>	3,097	6,2263
<i>qt_tok_dir_barra</i>	1,7573	2,9908
<i>sn_assoc</i>	2,4993	1,8551
<i>rank_alexa</i>	129,06	16430,6377
<i>rank_google</i>	2,0873	0,359

# Provas H3

Conjunto	Descrição	Conjunto	Descrição
A	Composto pelas 3 características baseadas em informações obtidas do DNS	B	Composto pelas 6 características denominadas especiais
C	Composto pelas 15 características léxicas mais comuns <sup>1</sup>	D	Composto por 32 características léxicas variáveis (não aparecem em todas as bases de dados)
E	Composto por todas as características léxicas	A+B+C	Composto pelos grupos A, B e C, totalizando 24 características
A+B+D	Composto pelos grupos A, B e D, totalizando 41 características	Todos	Todas as 56 características

# Provas H2 e H4

DMOZ/PhishTank		DMOZ/Blacklist	
Característica	InfoGain	Característica	InfoGain
<i>rank_google</i>	0.390339	<i>qt_tok_url_barra</i>	0.70944
<i>data_tempo_ativo</i>	0.348147	<i>data_tempo_ativo</i>	0.452587
<i>geo_localizacao</i>	0.228261	<i>qt_tok_dir_barra</i>	0.276703
<i>qt_tok_dom_ponto</i>	0.158864	<i>geo_localizacao</i>	0.2612
<i>qt_tok_url_barra</i>	0.14602	<i>qt_tok_dom_ponto</i>	0.237529
<i>qt_tok_dir_barra</i>	0.118791	<i>comp_arquivo</i>	0.233644
<i>comp_diretorio</i>	0.115648	<i>rank_alexa</i>	0.180241
<i>comp_url</i>	0.105754	<i>comp_dominio</i>	0.176853
<i>qt_tok_url_ponto</i>	0.096124	<i>ip_assoc</i>	0.162585
<i>rbf_check</i>	0.082557	<i>qt_tok_url_ponto</i>	0.133253
<i>comp_arquivo</i>	0.074864	<i>rank_google</i>	0.120826
<i>comp_dominio</i>	0.073538	<i>presenca_marca</i>	0.102384
<i>sn_assoc</i>	0.068434	<i>sn_assoc</i>	0.067489
<i>presenca_marca</i>	0.056488	<i>rbf_check</i>	0.067109
<i>rank_alexa</i>	0.051651	<i>comp_diretorio</i>	0.06236
<i>ip_assoc</i>	0.038537	<i>comp_url</i>	0.052208
<i>qt_tok_dir_ponto</i>	0.035866	<i>comp_params</i>	0.037848
<i>comp_params</i>	0.029835	<i>qt_tok_arq_ponto</i>	0.032119
<i>qt_tok_url_hifen</i>	0.029591	<i>qt_tok_dir_hifen</i>	0.03101
<i>qt_tok_dir_hifen</i>	0.028234	<i>qt_tok_url_igualdade</i>	0.02972
<i>qt_tok_url_ecomerc</i>	0.024787	<i>qt_tok_par_igualdade</i>	0.029292
<i>qt_tok_dom_hifen</i>	0.02455	<i>qt_tok_arq_til</i>	0.022689
<i>qt_tok_par_ecomerc</i>	0.022983	<i>qt_tok_url_hifen</i>	0.018402
<i>qt_params</i>	0.022983	<i>qt_tok_url_interrog</i>	0.014041