

# Phishing URL Detection Dataset: Feature Descriptions

This document contains detailed descriptions of the features present in the phishing URL detection dataset used to train a hybrid LSTM-CNN deep learning model.

---

## URL Structure Features

- **url**: The actual web address being analyzed.
  - **length\_url**: Total number of characters in the URL.
  - **length\_hostname**: Length of the domain name (hostname).
  - **ip**: Presence of an IP address instead of a domain name.
  - **nb\_dots**: Number of "." (dots) in the URL.
  - **nb\_hyphens**: Number of hyphens in the URL.
  - **nb\_at**: Number of "@" symbols.
  - **nb\_qm**: Number of "?" symbols.
  - **nb\_and**: Number of "&" symbols.
  - **nb\_or**: Number of "|" symbols.
  - **nb\_eq**: Number of "=" symbols.
  - **nb\_underscore**: Number of "\_" symbols.
  - **nb\_tilde**: Number of "~" symbols.
  - **nb\_percent**: Number of "%" symbols.
  - **nb\_slash**: Number of forward slashes "/".
  - **nb\_star**: Number of "\*" symbols.
  - **nb\_colon**: Number of ":" characters.
  - **nb\_comma**: Number of commas.
  - **nb\_semicolumn**: Number of semicolons.
  - **nb\_dollar**: Number of "\$" symbols.
  - **nb\_space**: Number of space characters.
  - **nb\_www**: Number of times "www" appears.
  - **nb\_com**: Number of times ".com" appears.
  - **nb\_dslash**: Number of "/" substrings.
  - **http\_in\_path**: Indicates if "http" appears in the path.
  - **https\_token**: Indicates if "https" is used incorrectly in the domain or path.
- 

## Ratio and Encodings

- **ratio\_digits\_url**: Ratio of digits in the entire URL.
  - **ratio\_digits\_host**: Ratio of digits in the hostname.
  - **punycode**: Checks if the domain uses punycode (internationalized domain names).
  - **port**: Presence of a port in the URL.
-

## Domain/Path Based Features

- **tld\_in\_path**: Top-Level Domain (TLD) appears in the path.
  - **tld\_in\_subdomain**: TLD appears in the subdomain.
  - **abnormal\_subdomain**: Irregular subdomain structure.
  - **nb\_subdomains**: Number of subdomains.
  - **prefix\_suffix**: Checks for prefixes/suffixes separated by "-" in domain.
  - **random\_domain**: Domain contains random strings.
  - **shortening\_service**: URL is shortened via services like bit.ly.
  - **path\_extension**: Path has suspicious extensions.
  - **nb\_redirection**: Number of redirections in URL.
  - **nb\_external\_redirection**: Number of redirections to external domains.
- 

## Word-Based Features

- **length\_words\_raw**: Total characters in all words in URL.
  - **char\_repeat**: Frequency of character repetition.
  - **shortest\_words\_raw**: Length of the shortest word in raw URL.
  - **shortest\_word\_host**: Shortest word in host.
  - **shortest\_word\_path**: Shortest word in path.
  - **longest\_words\_raw**: Length of the longest word in raw URL.
  - **longest\_word\_host**: Longest word in host.
  - **longest\_word\_path**: Longest word in path.
  - **avg\_words\_raw**: Average word length in full URL.
  - **avg\_word\_host**: Average word length in hostname.
  - **avg\_word\_path**: Average word length in path.
- 

## Suspicious Patterns

- **phish\_hints**: Presence of known phishing hints.
  - **domain\_in\_brand**: Brand name present in domain.
  - **brand\_in\_subdomain**: Brand present in subdomain.
  - **brand\_in\_path**: Brand present in path.
  - **suspicious\_tld**: TLDs often associated with phishing (e.g., .tk).
  - **statistical\_report**: Matches known phishing patterns from stats.
- 

## Hyperlink and Resource Analysis

- **nb\_hyperlinks**: Total number of hyperlinks.
- **ratio\_intHyperlinks**: Internal hyperlink ratio.
- **ratio\_extHyperlinks**: External hyperlink ratio.
- **ratio\_nullHyperlinks**: Ratio of empty/null hyperlinks.
- **nb\_extCSS**: Number of external CSS links.
- **ratio\_intRedirection**: Ratio of internal redirection.
- **ratio\_extRedirection**: Ratio of external redirection.
- **ratio\_intErrors**: Internal error links.

- **ratio\_extErrors**: External error links.
- 

## HTML and JavaScript Behavior

- **login\_form**: Presence of login forms.
  - **external\_favicon**: External favicon source.
  - **links\_in\_tags**: Hyperlink tags in HTML.
  - **submit\_email**: Form submits to email address.
  - **ratio\_intMedia**: Internal media links ratio.
  - **ratio\_extMedia**: External media links ratio.
  - **sfh**: Server Form Handler — suspicious values (e.g., blank).
  - **iframe**: Use of iframes.
  - **popup\_window**: Popup window detection.
  - **safe\_anchor**: Proportion of anchors with safe hrefs.
  - **onmouseover**: Use of onmouseover JavaScript.
  - **right\_click**: Right-click is disabled.
  - **empty\_title**: Page has empty title.
  - **domain\_in\_title**: Domain appears in title.
- 

## Domain and Reputation

- **domain\_with\_copyright**: Domain has a copyright symbol.
  - **whois\_registered\_domain**: WHOIS lookup succeeded.
  - **domain\_registration\_length**: Time until domain expires.
  - **domain\_age**: Domain age in days.
  - **web\_traffic**: Site's web traffic rank.
  - **dns\_record**: Presence of DNS records.
  - **google\_index**: Whether site is indexed by Google.
  - **page\_rank**: Page rank of site.
- 

## Labels

- **status**: Class label — `legitimate` or `phishing`.