

Nithin S
221IT085

IT352 Course Project Module 1

Verifying existing URLs with VirusTotal and Extracting Features from URL Dataset to build a new dataset

Phishing Site URLs: Dataset which contains Phishing urls and non phishing urls.

Total URLs: 549362

Total Unique URLs: 507195

URL	Label
unique URLs	good and bad URLs classes
507195 unique values	good 72% bad 28%
nobell.it/70ffb52d079109dca5664cce6f317373782/login.SkyPe.com/en/cgi-bin/verification/login/70ffb52d...	bad
www.dghjdgf.com/paypal.co.uk/cycgi-bin/webcmd=_home-customer&nav=1/loading.php	bad
serviciosbys.com/paypal.cgi.bin.get-into.herf.secure.dispatch35463256r321654641dsf654321874/herf/h...	bad

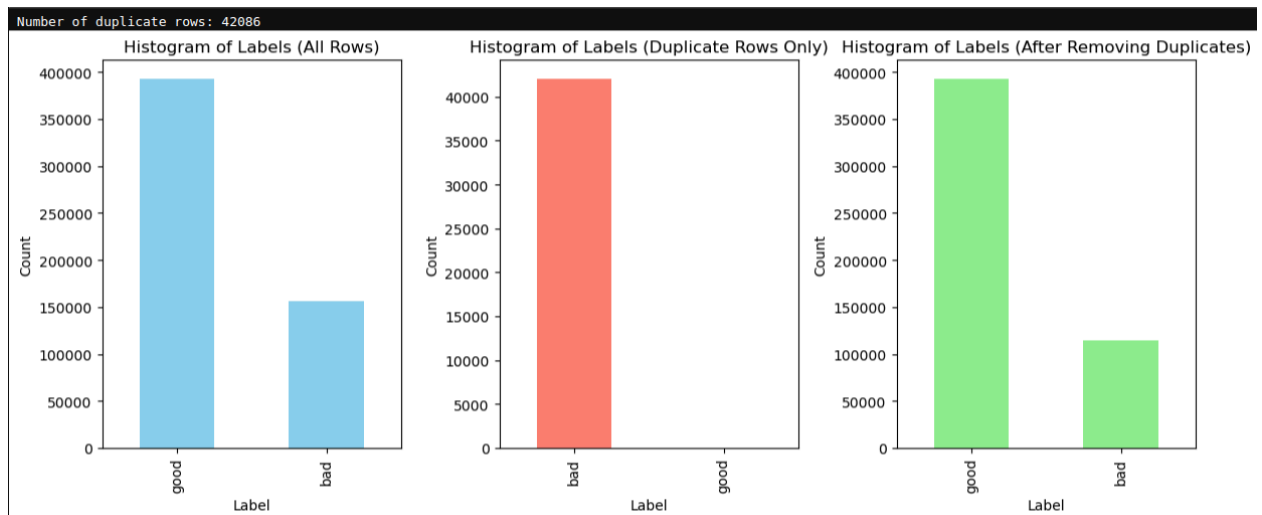
Histogram about the classes in the URLs Dataset

Total Rows present in given dataset: 549346

No of duplicates : 42086

Percentage of duplicates: 7.66 %

Total Unique URLs: 507195



No Null Values were found in the given URLs Dataset

```
phish_data.isnull().sum() # there is no missing values
```

```
URL      0
Label    0
dtype: int64
```

Python Script to verify the labelling of the dataset

Please make sure to add required API Keys before running the code.

Virus Total API : <https://docs.virustotal.com/reference/scan-url>

Virus Total API Endpoint has strict rate limits

API QUOTA ALLOWANCES FOR YOUR USER

You own a standard free end-user account. It is not tied to any corporate group and so it does not have access to Premium services. You are subjected to the following limitations:

Access level	△ Limited , standard free public API	Upgrade to premium
Usage	Must not be used in business workflows, commercial products or services.	
Request rate	4 lookups / min	
Daily quota	500 lookups / day	
Monthly quota	15.5 K lookups / month	

So, we will verify random 500 urls from the dataset without repetition and update them. We will have a **time delay of 60s** between each URL Verification request.

CODE

```
def check_url_virustotal(url, default_label):
    endpoint = "https://www.virustotal.com/api/v3/urls"
    headers = {"x-apikey": os.getenv("VIRUS_TOTAL_API_KEY")}
    try:
        response = requests.post(endpoint, headers=headers, data={"url":
url})
        if response.status_code != 200:
            return default_label
        analysis_id = response.json()["data"]["id"]
        result_endpoint =
f"https://www.virustotal.com/api/v3/analyses/{analysis_id}"
        result_response = requests.get(result_endpoint, headers=headers)
        if result_response.status_code != 200:
            return default_label
```

```

        result_data =
result_response.json()["data"]["attributes"]["results"]
        malicious_count = sum(1 for scan in result_data.values() if
scan["category"] == "malicious")
        return "bad" if malicious_count > 0 else "good"
    except Exception:
        return default_label

```

Output

```

(env) nithin@pavilion:~/Codes/Projects/Phishing Website Detection$ python virustotal.py
[1/500] Checked URL: vehiculapress.com/montreal/writers/richler.html
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:19:04

[2/500] Checked URL: hfmacabados.com/verify-boainfps/bank-account/login/info.php
Original Label: bad -> Updated Label: good | Not Match
Analysis Timestamp: 2025-03-05 00:20:07

[3/500] Checked URL: pix.eosphotonan.com/
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:21:09

[4/500] Checked URL: 'www.gswzmb.com/down/js/7us.battle.net/login/en/?ref=http%3A%2F%2Fus.battle.net%2Fd3%2Fen%2Findex&app=com-d3'
Original Label: bad -> Updated Label: bad | Match
Analysis Timestamp: 2025-03-05 00:22:11

[5/500] Checked URL: torosakiskan.com/modules/telekom/telekom_deutschland_gnbh
Original Label: bad -> Updated Label: good | Not Match
Analysis Timestamp: 2025-03-05 00:23:13

[6/500] Checked URL: btjunkie.org/torrent/Stolen-Babies-There-Be-Squabbles-Ahead-2006/37819f37750860e9ab9871e90ffbcc77de086c07ab56
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:24:15

[7/500] Checked URL: en.wikipedia.org/wiki/Getchell_Mine
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:25:18

[8/500] Checked URL: ottenheim.cdlex.org/
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:26:20

[9/500] Checked URL: www.bus.wisc.edu/ASRMI/
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:27:23

[10/500] Checked URL: illicoweb.videotron.com/illicoweb/channels/Mlle
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:28:26

[11/500] Checked URL: litigation-essentials.lexisnexis.com/webcd/app?action=DocumentDisplay&crawlid=1&doctype=cite&docid=55+S.+Cal.+L.+Rev.+895&srctype=snl&srcid=3B15&key=a8575485bbb18628e7e
254c8aca6bdeb
Original Label: good -> Updated Label: good | Match
Analysis Timestamp: 2025-03-05 00:29:28

```

Total mismatches : 120 out of 500

Percentage of mismatches: $120/500 * 100 = 24\%$

Due to size limitations I have uploaded the datasets in my google drive.

Extracted Features Dataset Google Drive Link:

https://drive.google.com/file/d/119L_eJDb8Oizm4jv3NM0tfW-Um-7sSDC/view?usp=sharing

Z Score Normalized Dataset Google Drive Link:

https://drive.google.com/file/d/1P8sOZi6j7JeM_LNLFmChU9FPHnyo_u1X/view?usp=sharing

List of Features planned to be extracted: 116 Features - (Without dropping repeated columns)

Full URL Length: Total number of characters in the entire URL.

Hostname (Domain) Length: Total number of characters in the domain name part.

Directory Length: Number of characters in the folder or path part of the URL.

File Name Length: Number of characters in the file name portion (if any).

Parameters Length: Number of characters in the query string (everything after "?").

TLD Length: Length (in characters) of the top-level domain (for example, "com" or "org").

Dot ('.') Count: Number of periods used (often separating subdomains or domain parts).

Hyphen ('-') Count: Number of hyphens used.

Underscore ('_') Count: Number of underscores.

Slash ('/') Count: Number of forward slashes.

Question Mark ('?') Count: How many "?" appear.

Equal Sign ('=') Count: Number of equal signs.

At Sign ('@') Count: Number of "@" symbols.

Ampersand ('&') Count: How many "&" symbols appear.

Exclamation Mark ('!') Count: Number of "!" symbols.

Space Count: Number of space characters.

Tilde ('~') Count: Number of tilde characters.

Comma (',') Count: How many commas appear.

Plus Sign ('+') Count: Number of "+" symbols.

Asterisk ('*') Count: Number of "*" symbols.

Hashtag ('#') Count: Number of "#" symbols.

Dollar Sign ('\$') Count: Number of "\$" symbols.

Percent Sign ('%') Count: Number of "%" symbols.

Common Terms Occurrence: Counts for terms such as "www", ".com", "http", and "/" that usually appear only once in normal URLs.

Email in URL: A flag indicating if an email address is embedded in the URL.

HTTPS Token: Checks if the URL uses "https" (a sign of secure connections).

IP Address in URL: A binary check to see if an IP address is used instead of a domain name.

Punycode Usage: Checks whether the domain uses punycode (which can mask its true characters).

Port Number Presence: A flag indicating if the URL explicitly shows a port (like ":80" or ":443").

TLD Position: Verifies that the top-level domain is in the right place (it should not appear in the wrong section like the path or subdomain).

Abnormal Subdomains: Detects unusual subdomain patterns (for example, variations of “www” that include numbers).

Number of Subdomains: Counts how many subdomains are present.

Prefix/Suffix with Hyphen: Checks if the domain uses hyphens to separate extra words (which might be used to mimic legitimate sites).

Random Domain Indicator: Determines if the domain seems to be made up of random characters.

URL Shortening Service: A flag to see if a URL shortener (like bit.ly) is used, which can hide the true destination.

Path Extension Check: Looks for suspicious file extensions (such as “.exe” or “.js”) in the URL path.

Suspicious TLD: Checks if the top-level domain is among those known to be risky.

Digit Ratio in Full URL: Proportion of digit characters compared to the total characters in the URL.

Digit Ratio in Hostname: Proportion of digits in the domain name itself.

Word Count: Number of words found in the full URL, the hostname, or the path.

Shortest & Longest Word: Identification of the shortest and longest word in the URL parts.

Average Word Length: The average length of words in the URL, hostname, or path.

Phish Hints: Counts occurrences of suspicious or phishing-related keywords (like “login”, “admin”, “signin”, etc.).

Brand Names in URL:

In the Domain: Presence of well-known brand names can be a sign of legitimacy.

In the Subdomain or Path: Their appearance here may indicate an attempt to deceive.

Domain in Page Title/Copyright: Checks if the domain name appears in the webpage title or copyright text (a sign of legitimacy).

Redirection Count: Total number of times the URL redirects to another page.

External Redirections: How many of these redirects go to a different domain.

Internal vs. External Hyperlinks Ratio: Compares links that point within the same site to those that point to external sites.

Null Hyperlinks Ratio: Proportion of links that lead nowhere (empty links).

Media Links Ratio: Ratio of media (images, videos, etc.) hosted on the same domain versus externally.

Connection Errors Ratio: Ratio of hyperlinks that result in errors (broken links).

Number of Hyperlinks: Total links present on the webpage.

External CSS Files Count: Number of CSS files linked from outside the domain.

Login Forms Presence: Checks for login forms, especially those with empty or suspicious action attributes.

External Favicon: Whether the page uses a favicon (the small icon in the browser tab) from an external source.

Invisible iFrame: Detects hidden iframe elements that might load content from another domain.

Pop-up Windows: Looks for pop-up windows that include text fields (which can be a sign of phishing).

Unsafe Anchors: Counts anchor (<a>) tags that use unsafe links (e.g., "javascript:" or "#").

Right-Click Blocking: Checks for scripts that disable the right-click function (which can hide page source).

Empty Title: Flags if the webpage has no title tag.

WHOIS Registration: Whether the domain is found in the WHOIS database (a missing record is a red flag).

Domain Registration Length: The number of years for which the domain is registered (short registration periods can be suspicious).

Domain Age: How long the domain has been active.

DNS Record Check: Verifies that the domain has proper DNS records.

Google Index: Checks if the URL or domain is indexed by Google (phishing sites are often not).

Page Rank: An estimate of the webpage's popularity.

Web Traffic: An indicator (like Alexa ranking) showing the number of visitors.

Additionally, one study mentions a "statistical report" feature that checks if the domain's IP matches known top phishing domains.

Vowel Count in Domain: Number of vowels in the domain name.

Domain in IP Format: Whether the domain is written as an IP address.

"Server" or "Client" in Domain: Checks if these words appear in the domain name, which can hint at its purpose.

Domain Lookup Response Time: How long it takes to get a response when looking up the domain.

SPF Record: Checks if the domain has an SPF record (helps validate email sources).

ASN (Autonomous System Number): A number that identifies the network the domain's IP belongs to.

Domain Activation Time: How many days have passed since the domain was first activated.

Domain Expiration Time: How many days remain until the domain expires.

Number of Resolved IPs: How many IP addresses are returned when the domain is looked up.

Nameservers Count: Number of DNS nameservers linked to the domain.

MX Servers Count: Number of mail servers associated with the domain.

TTL (Time-To-Live) of Hostname: The DNS record's lifetime.

Valid TLS/SSL Certificate: Whether the site has a proper secure certificate.

URL Shortened Flag: Whether the URL has been shortened (also noted earlier under security).

TLD Present in Parameters: Checks if a top-level domain appears within the URL parameters (which is unusual).

Number of Parameters: Count of key-value pairs or parameters present in the URL query string.

Check for duplicate Columns

```
Column Index: 14, Column Name: tilde_count  
Column Index: 28, Column Name: https_token  
Column Index: 60, Column Name: brand_in_subdomain  
Column Index: 86, Column Name: whois_registration  
Column Index: 87, Column Name: domain_registration_length  
Column Index: 88, Column Name: domain_age  
Column Index: 96, Column Name: server_or_client_in_domain  
Column Index: 98, Column Name: asn  
Column Index: 99, Column Name: domain_activation_time  
Column Index: 100, Column Name: domain_expiration_time
```

Check for duplicate rows

```
5339  
9567  
16871  
18507  
18965  
18967  
19005  
19398  
19673  
19897  
22329  
22413  
22605  
22755  
22762  
22824  
22892
```

Normalization

I have normalized all columns , since all were numerical.


```
Standardization (Z-score normalization) applied to: ['full_url_length', 'hostname_length', 'ip_address_in_url', 'dot_count', 'hyphen_count', 'underscore_count', 'slash_count', 'question_mark_count', 'equal_count', 'at_count', 'ampersand_count', 'exclamation_count', 'space_count', 'tilde_count', 'comma_count', 'plus_count', 'asterisk_count', 'hashtag_count', 'dollar_count', 'percent_count', 'vertical_bar_count', 'colon_count', 'semicolon_count', 'www_occurrence', 'com_occurrence', 'http_occurrence', 'double_slash_occurrence', 'https_token', 'digit_ratio_full_url', 'digit_ratio_hostname', 'punycode_usage', 'port_number_presence', 'tld_in_path', 'tld_in_hostname', 'abnormal_subdomains', 'number_of_subdomains', 'prefix_suffix_hyphen', 'random_domain_indicator', 'url_shortening_service', 'path_extension_check', 'redirection_count', 'external_redirection_count', 'word_count_url', 'word_count_hostname', 'word_count_path', 'char_repeat_url', 'char_repeat_hostname', 'char_repeat_path', 'shortest_word_url', 'shortest_word_hostname', 'shortest_word_path', 'longest_word_url', 'longest_word_hostname', 'longest_word_path', 'average_word_length_url', 'average_word_length_hostname', 'average_word_length_path', 'phish_hints', 'brand_in_domain', 'brand_in_subdomain', 'brand_in_path', 'suspicious_tld', 'statistical_report', 'number_of_hyperlinks', 'null_hyperlinks_ratio', 'external_css_files_count', 'internal_redirection_ratio', 'external_redirection_ratio', 'internal_errors_ratio', 'external_errors_ratio', 'login_forms_presence', 'external_favicon', 'internal_link_tags_ratio', 'submit_to_email', 'internal_hyperlink_ratio', 'external_hyperlink_ratio', 'internal_media_ratio', 'external_media_ratio', 'sfh', 'invisible_iframe', 'pop_up_windows', 'unsafe_anchors', 'right_click_blocking', 'empty_title', 'domain_in_copyright', 'whois_registration', 'domain_registration_length', 'domain_age', 'directory_length', 'file_name_length', 'parameters_length', 'tld_length', 'email_in_url', 'vowel_count_in_domain', 'domain_in_ip_format', 'server_or_client_in_domain', 'domain_lookup_response_time', 'asn', 'domain_activation_time', 'domain_expiration_time', 'number_of_resolved_ips', 'nameservers_count', 'ttl_hostname', 'ttl_hostname', 'tls_ssl_certificate', 'tld_present_in_parameters', 'number_of_parameters', 'dns_record_check', 'media_links_ratio', 'connection_errors_ratio', 'mx_servers_count', 'spf_record', 'domain_in_title', 'web_traffic', 'google_index', 'page_rank']
```

Handle Missing Values

```
df.isnull().sum()

URL                0
full_url_length    0
hostname_length    0
ip_address_in_url  0
dot_count          0
..
domain_in_title    0
web_traffic        0
google_index       0
page_rank          0
Label             0
Length: 117, dtype: int64
```

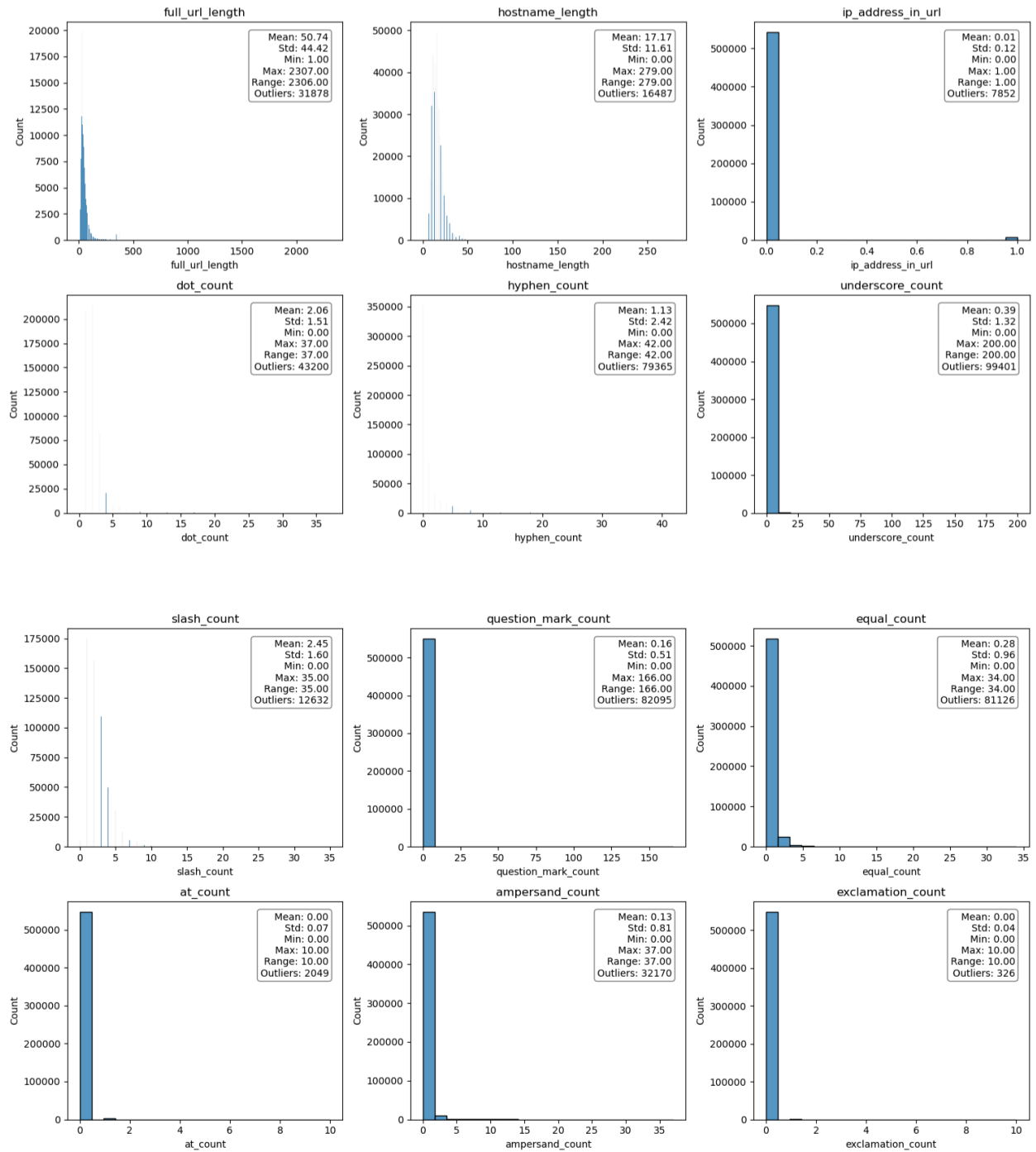
NO missing values were found

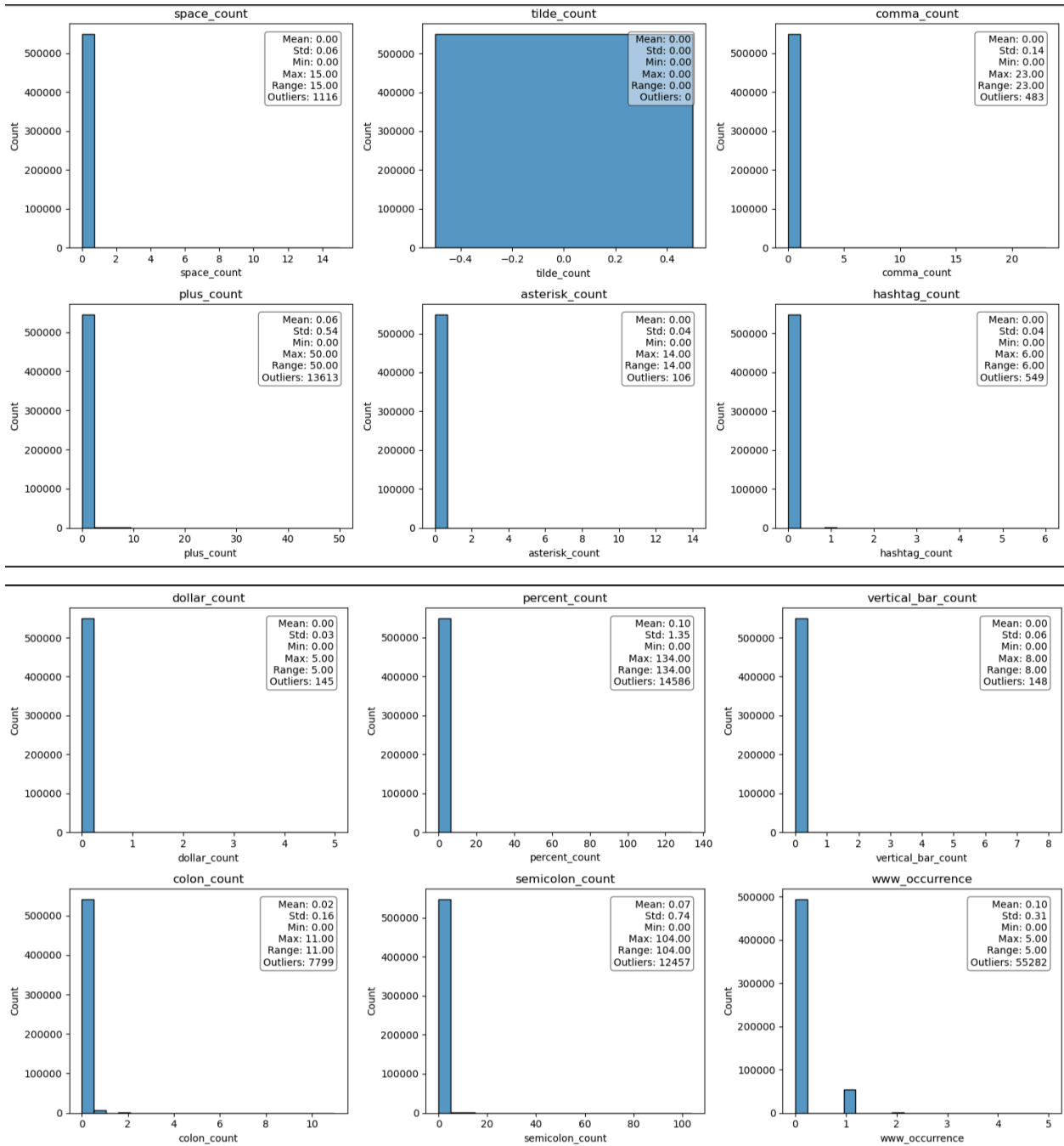
Correlation Map

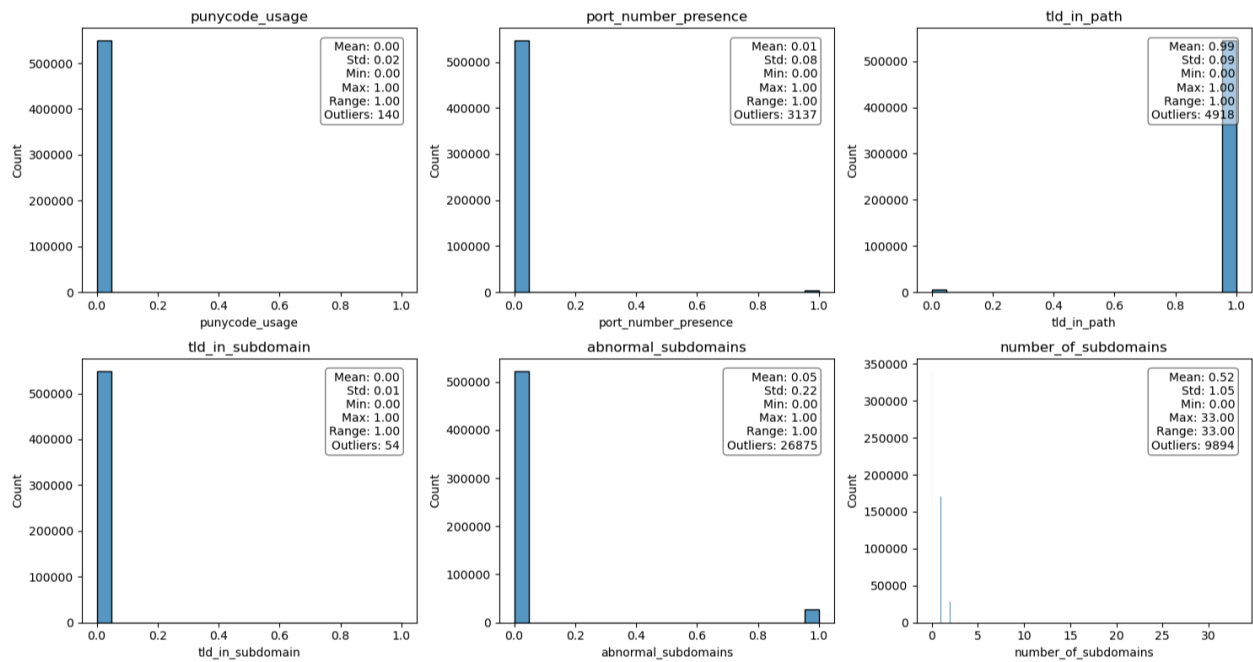
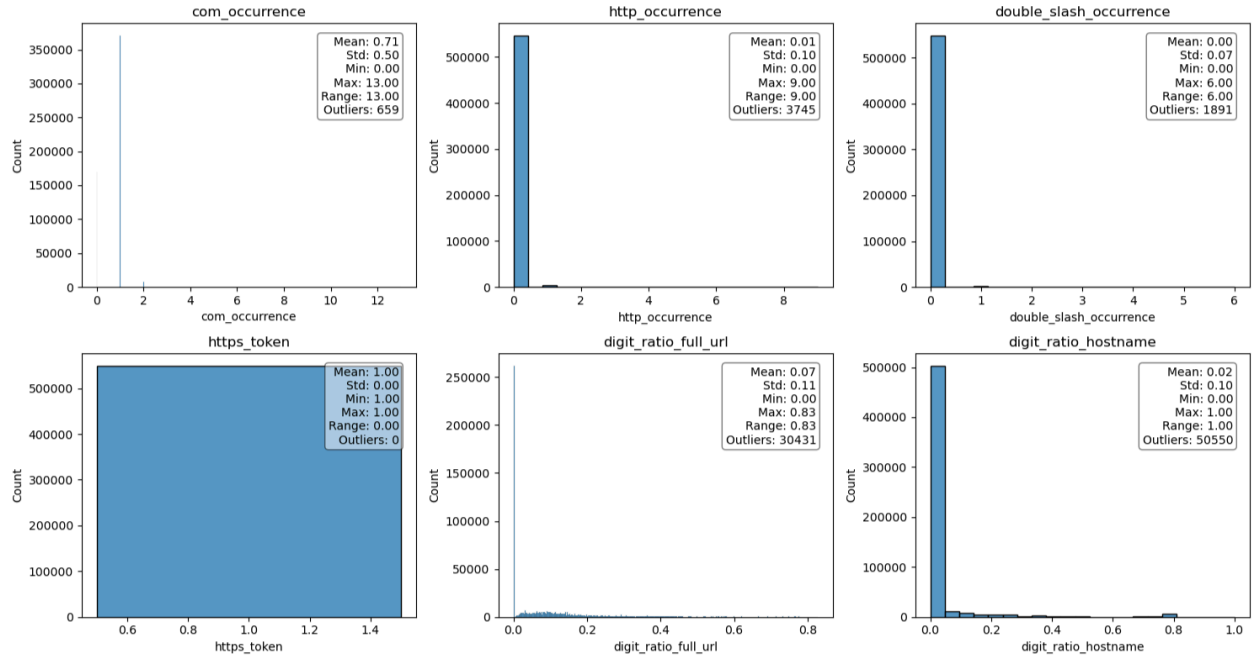


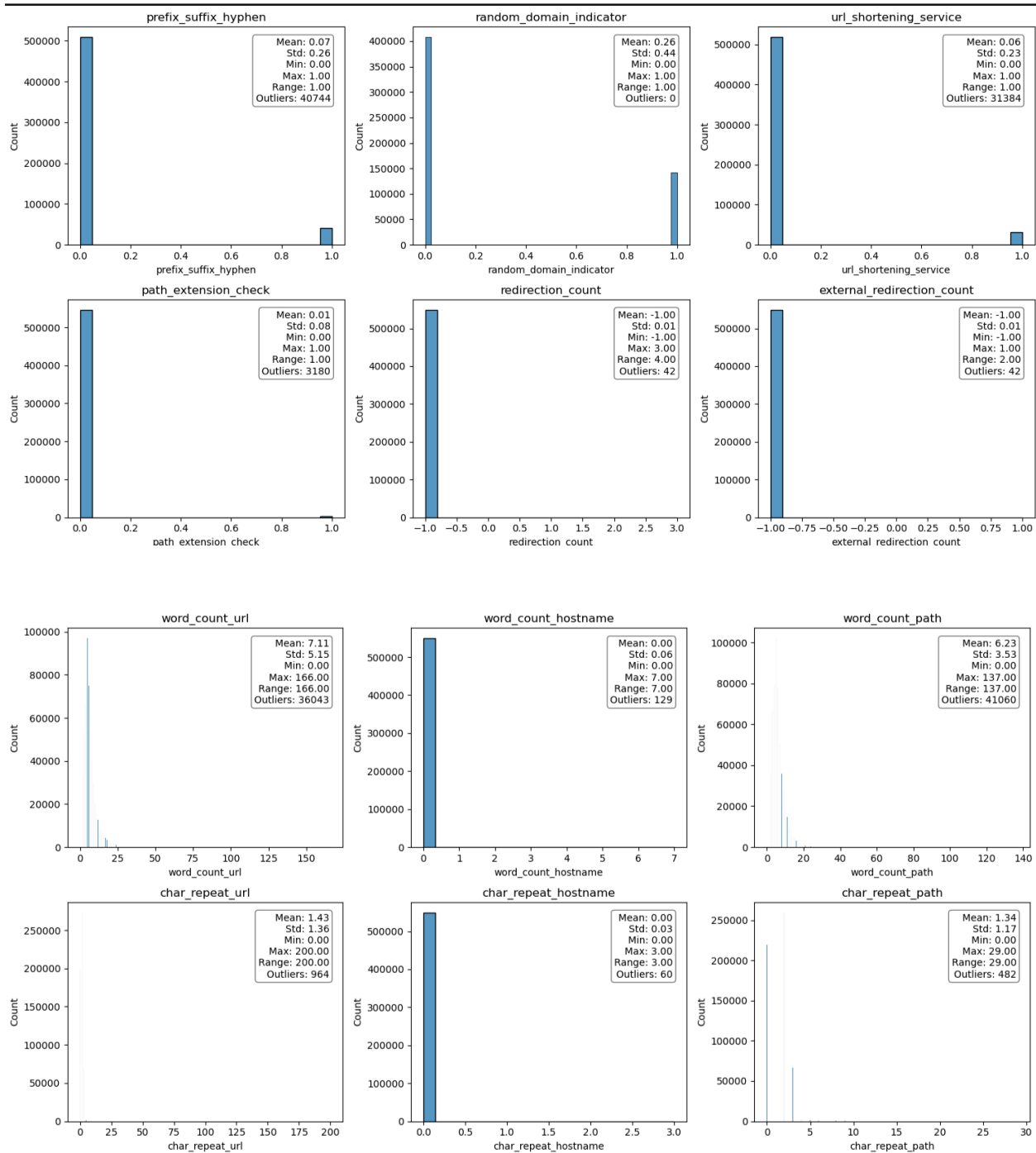
Please download the jpeg image from moodle and zoom into it for better view.

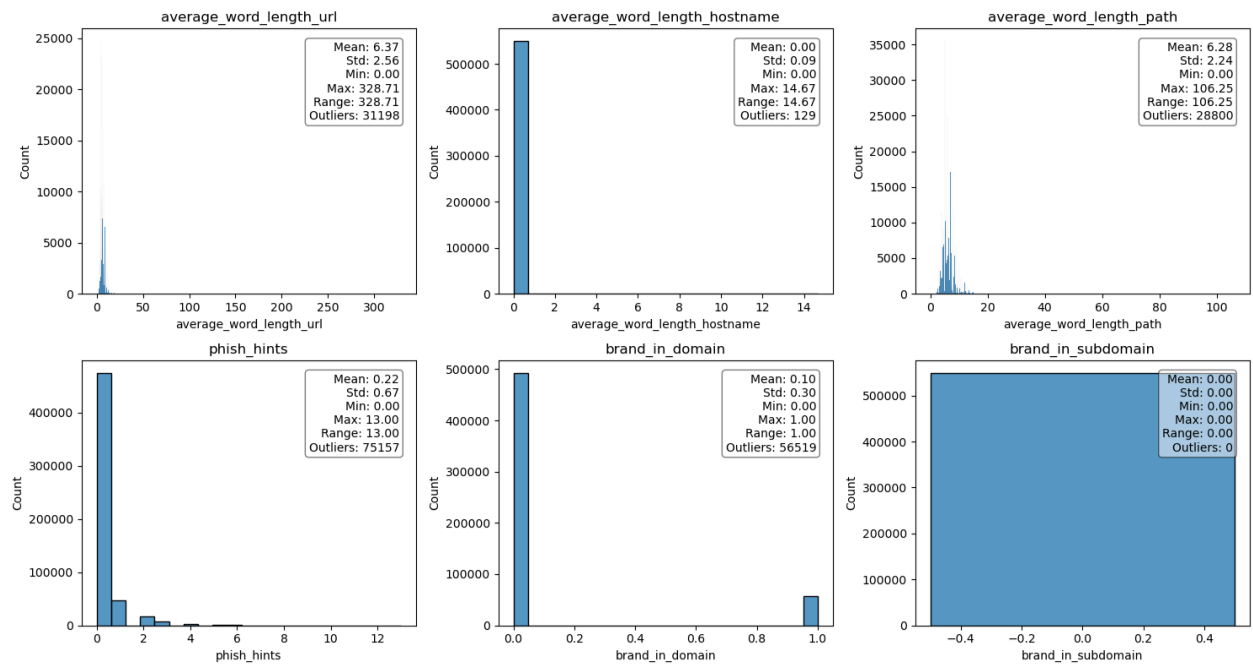
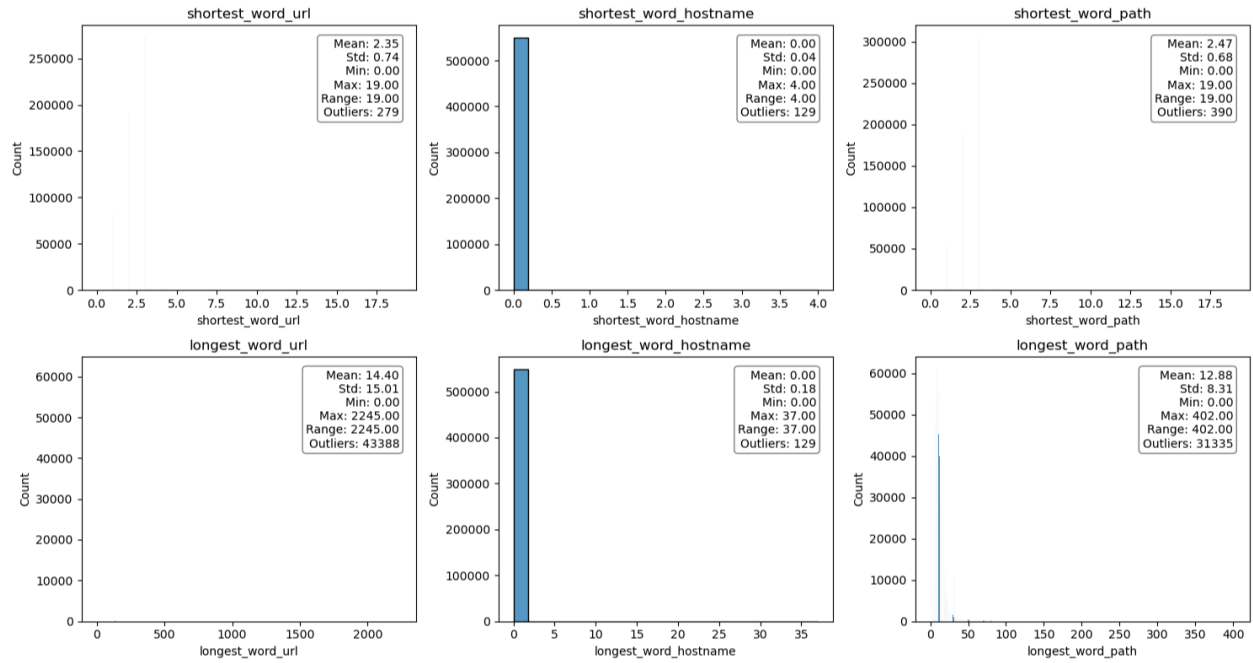
Distribution of Each Feature with its mean, standard deviation, min, max, range and outlier count

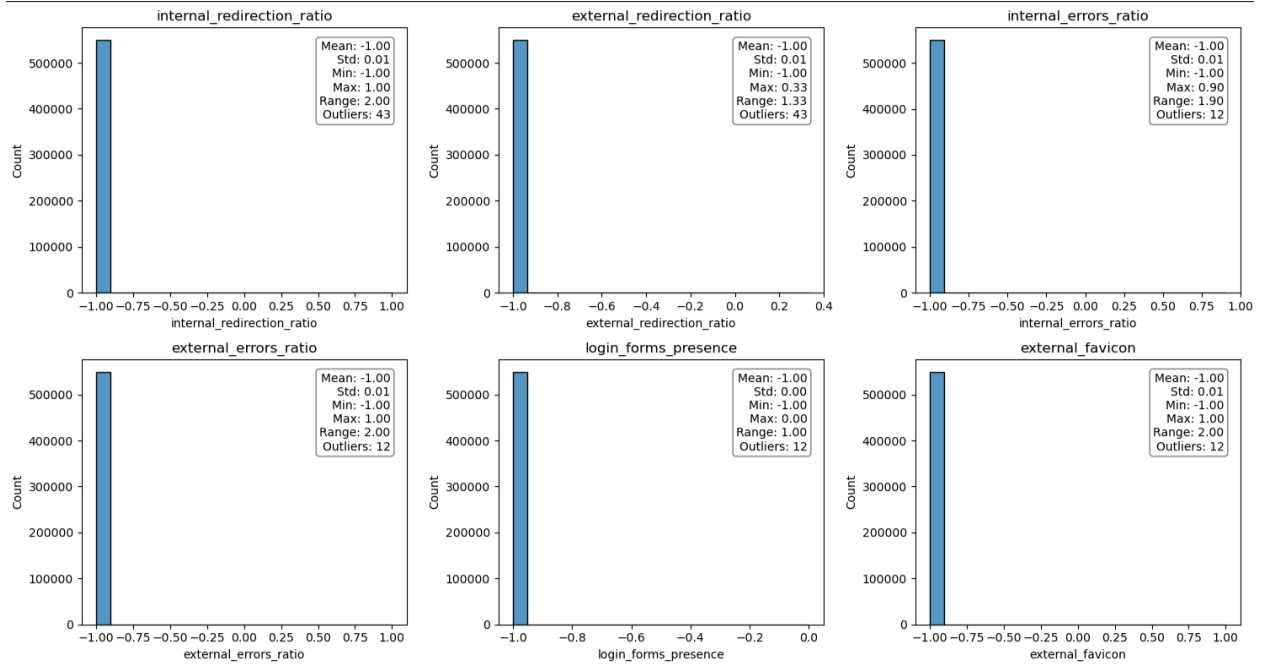
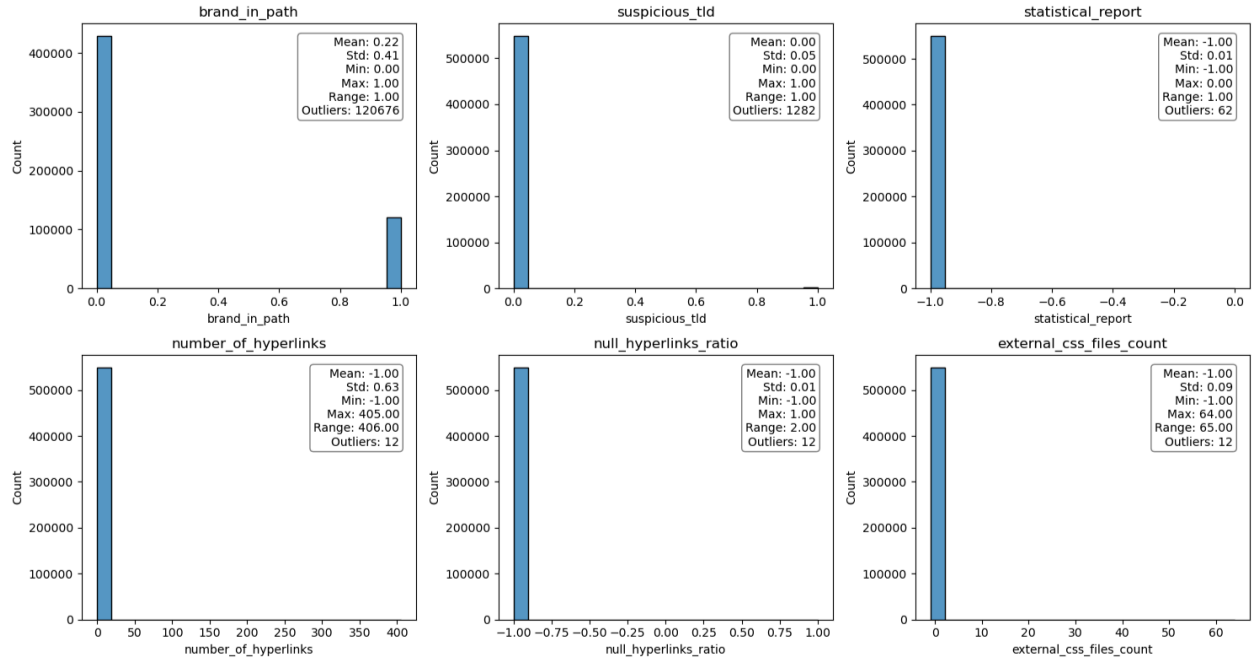


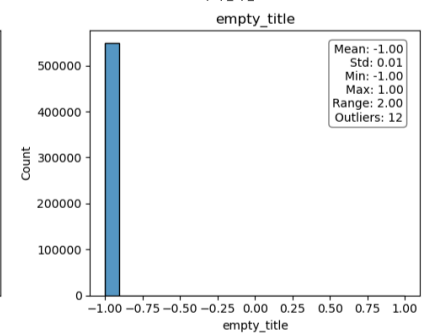
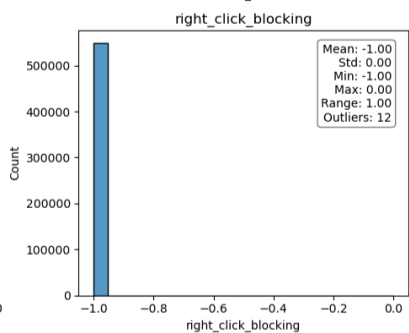
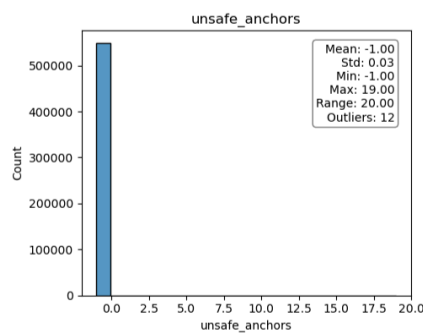
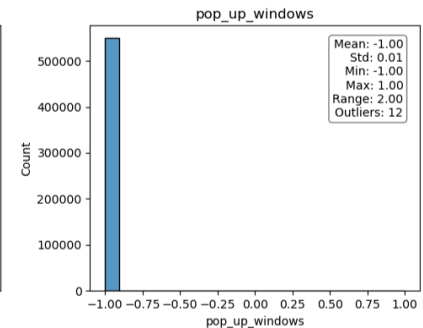
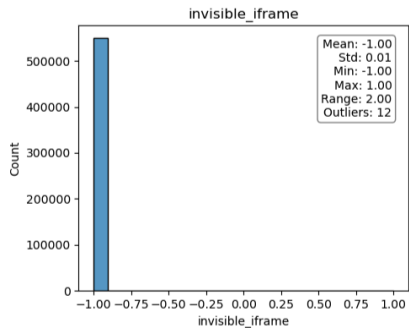
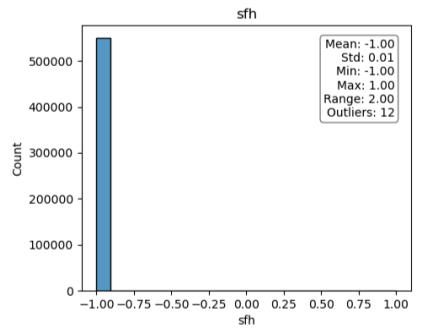
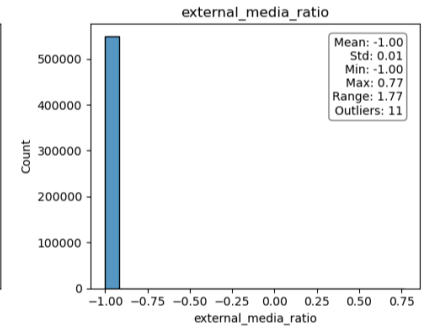
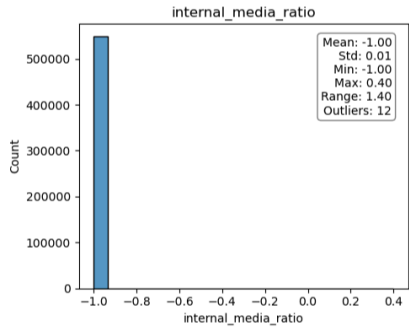
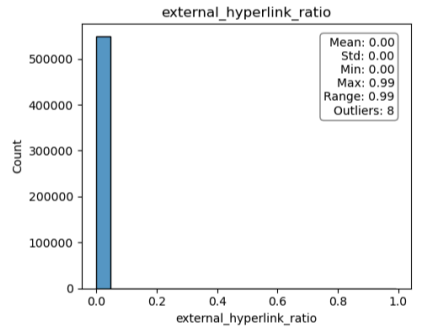
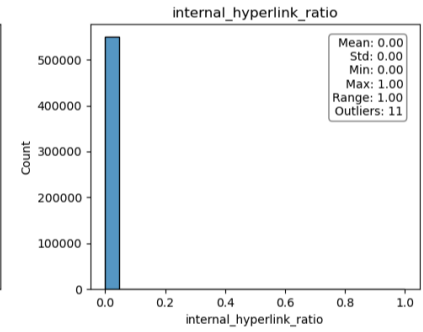
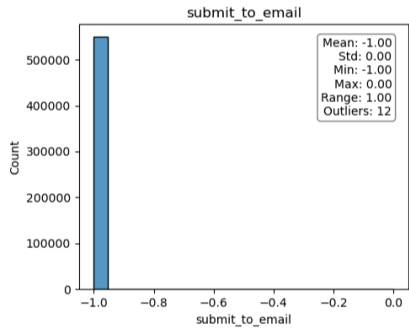
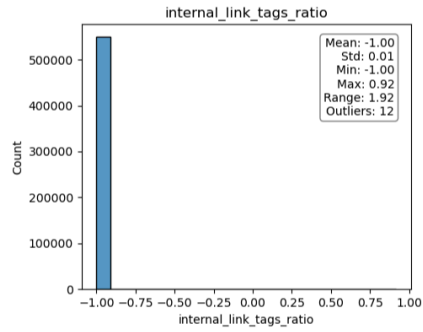


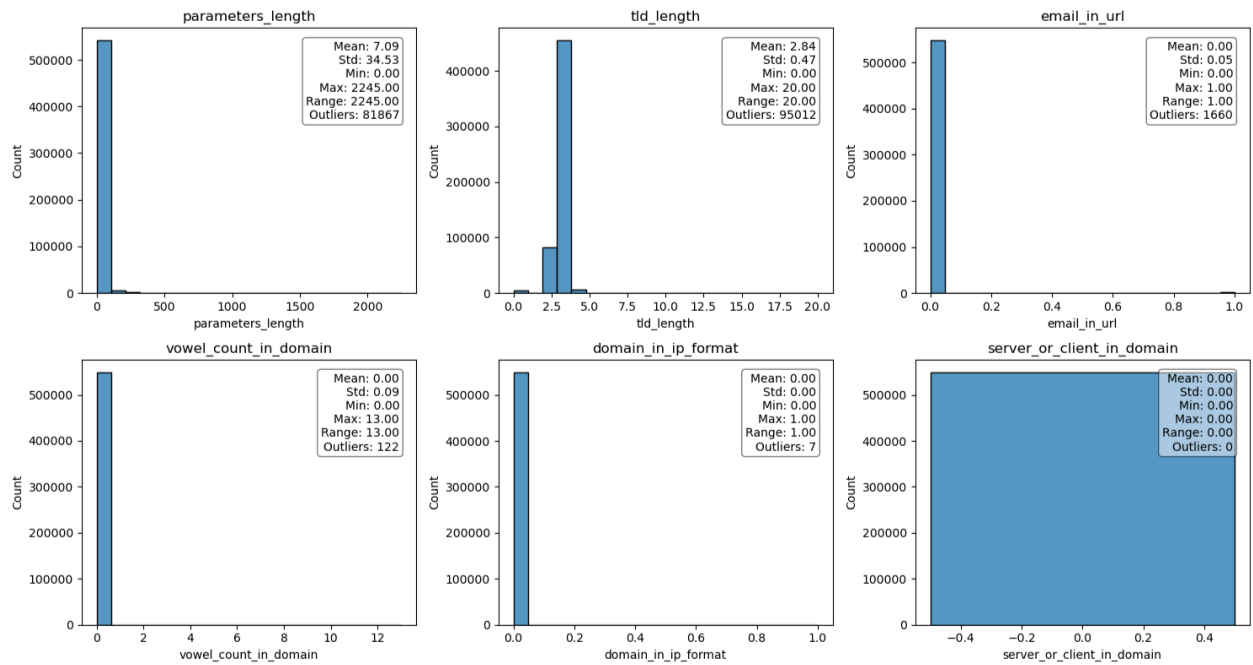
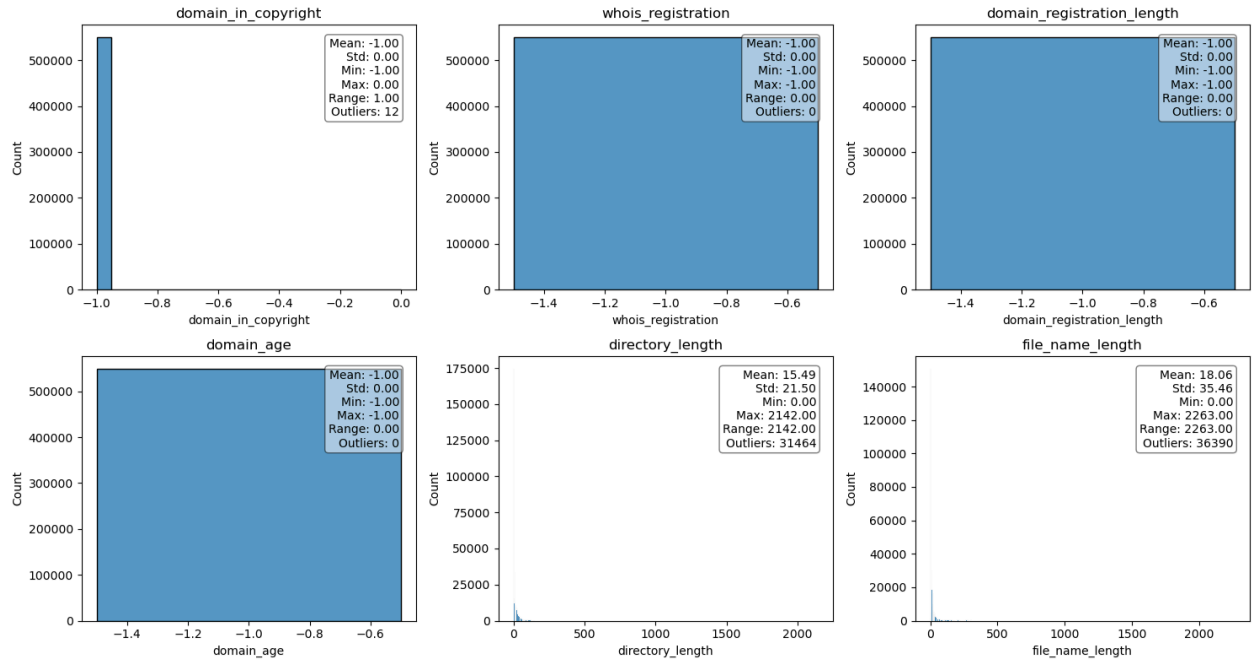


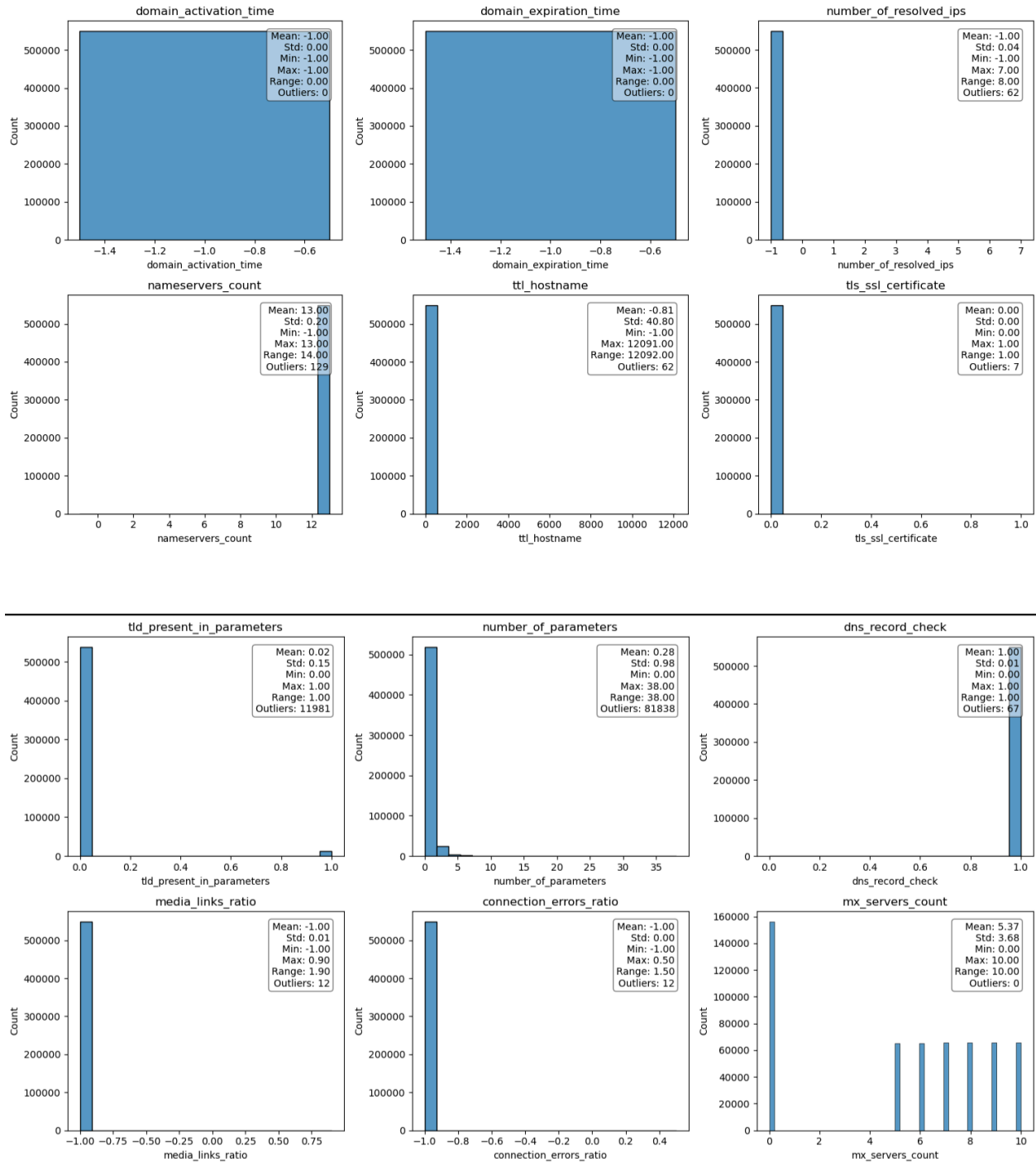






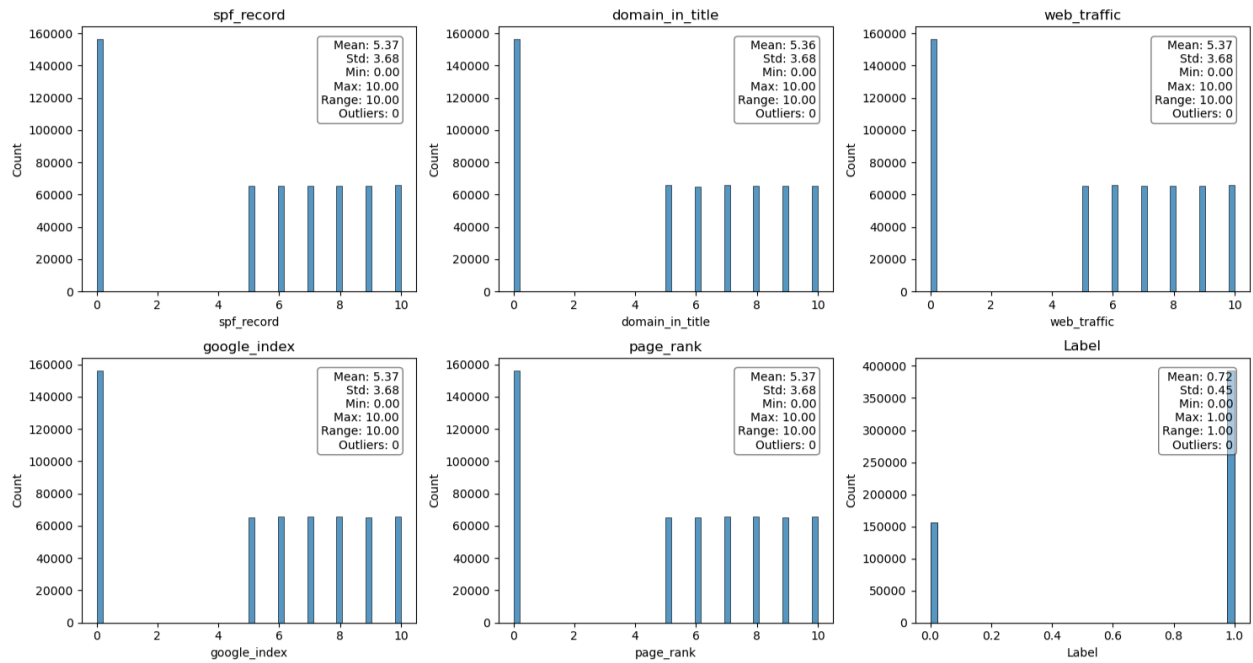






We will z score normalize only those features which dont have gaussian distribution and have a huge range.

We will drop duplicate columns and highly correlated features.



Drop Repeated Columns

```
Column Index: 14, Column Name: tilde_count
Column Index: 28, Column Name: https_token
Column Index: 60, Column Name: brand_in_subdomain
Column Index: 86, Column Name: whois_registration
Column Index: 87, Column Name: domain_registration_length
Column Index: 88, Column Name: domain_age
Column Index: 96, Column Name: server_or_client_in_domain
Column Index: 98, Column Name: asn
Column Index: 99, Column Name: domain_activation_time
Column Index: 100, Column Name: domain_expiration_time
```

Drop highly correlated Columns with $\text{corr} > 0.9$

connection_errors_ratio,
'internal_link_tags_ratio',
'sfh',
'Nameservers_count',
'Pop_up_windows',
'internal_redirection_ratio',
'External_favicon',
'Internal_media_ratio',
'External_errors_ratio',
'External_redirection_count',
'dns_record_check',
'right_click_blocking',
'External_redirection_ratio',
'internal_errors_ratio',
'Domain_in_copyright',
'Average_word_length_hostname',
'number_of_parameters',
'Vowel_count_in_domain',
'unsafe_anchors',
'Media_links_ratio',
'login_forms_presence',
'Empty_title',
'Invisible_iframe',
'Submit_to_email',
'longest_word_hostname',
'external_media_ratio'

