

Business Problem Understanding Document

1. Business Problem Summary

Phishing websites are fraudulent sites that attempt to steal sensitive information such as usernames, passwords, and financial details by masquerading as legitimate websites. The ability to detect phishing websites accurately is crucial for cybersecurity, as phishing attacks lead to financial losses, identity theft, and reputational damage for individuals and organizations.

The goal of this project is to develop a **phishing website detection model** using data-driven techniques. By analyzing various website features, the model will help identify phishing attempts and enhance online security. This solution is particularly valuable for businesses, financial institutions, and security firms seeking to protect users from cyber threats.

2. Key Insights from Literature on Phishing

- **Common Characteristics of Phishing Websites:**
 - Use of misleading domain names (e.g., typosquatting, similar-looking URLs).
 - Excessive redirections and shortened URLs.
 - Presence of deceptive pop-ups and fake login forms.
 - Lack of HTTPS security and valid SSL certificates.
 - **Detection Challenges:**
 - Phishing tactics evolve rapidly, making rule-based detection less effective.
 - Traditional blacklist-based approaches fail to detect new phishing sites.
 - Feature extraction from website content, URL structure, and HTML is complex.
 - **Potential Solutions:**
 - **Machine Learning-Based Detection:** Utilizing classifiers trained on website attributes.
 - **Heuristic Analysis:** Identifying suspicious patterns in URLs, domain registration, and HTTPS usage.
 - **Real-Time Monitoring:** Combining AI with real-time scanning to detect zero-day phishing sites.
-

Dataset Exploration Report

1. Overview of the Dataset

The dataset contains various features extracted from website URLs, HTML structure, and security indicators. The primary objective is to use these features to classify websites as either **phishing (malicious) or legitimate (safe)**.

- **Number of Features:** [To be determined from dataset analysis]
- **Types of Data:**
 - **Numerical Features:** Metrics such as domain age, URL length, number of subdomains.
 - **Categorical Features:** Presence of HTTPS, presence of login forms.
 - **Binary Features:** Indicators like "Has IP address in URL" (Yes/No).
- **Target Variable Distribution:** The dataset contains labeled data indicating whether a website is phishing (1) or legitimate (0). A distribution analysis will help identify any class imbalances that could affect model performance.

2. Description of Individual Features

Each feature contributes to phishing detection by identifying suspicious behaviors or security flaws. Key features include:

- **URL-Based Features:**
 - Length of URL, presence of special characters, number of subdomains.
 - Whether the domain uses an IP address instead of a hostname (common in phishing sites).
- **Domain and Hosting Information:**
 - Age of the domain (newly registered domains are often phishing sites).
 - WHOIS information consistency and SSL certificate validity.
- **Website Content & HTML Features:**
 - Presence of deceptive elements like fake login forms.
 - Frequency of redirects and presence of pop-ups.
- **Network and Security Indicators:**
 - Use of HTTPS (secure connection vs. insecure HTTP).
 - Presence in known phishing blacklists.

Conclusion

This dataset exploration helps establish a foundation for phishing detection by identifying important patterns in website attributes.

--- Internship project (phishing data set) ---

```
[22]: import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
from scipy import stats
from sklearn.preprocessing import LabelEncoder, OneHotEncoder, StandardScaler, MinMaxScaler
df=pd.read_csv('dataset_phishing.csv')
```

```
[2]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11430 entries, 0 to 11429
Data columns (total 89 columns):
#   Column                Non-Null Count  Dtype
---  -
0   url                    11430 non-null  object
1   length_url             11430 non-null  int64
2   length_hostname        11430 non-null  int64
3   ip                     11430 non-null  int64
4   nb_dots                 11430 non-null  int64
5   nb_hyphens              11430 non-null  int64
6   nb_at                   11430 non-null  int64
```

```

17:  def __init__(self):
18:      self.__dict__ = {}
19:      self.__dict__['name'] = 'None'
20:      self.__dict__['age'] = 'None'
21:      self.__dict__['sex'] = 'None'
22:      self.__dict__['height'] = 'None'
23:      self.__dict__['weight'] = 'None'
24:      self.__dict__['blood'] = 'None'
25:      self.__dict__['marriage'] = 'None'
26:      self.__dict__['children'] = 'None'
27:      self.__dict__['education'] = 'None'
28:      self.__dict__['income'] = 'None'
29:      self.__dict__['work'] = 'None'
30:      self.__dict__['religion'] = 'None'
31:      self.__dict__['race'] = 'None'
32:      self.__dict__['ethnicity'] = 'None'
33:      self.__dict__['nationality'] = 'None'
34:      self.__dict__['citizenship'] = 'None'
35:      self.__dict__['residence'] = 'None'
36:      self.__dict__['employment'] = 'None'
37:      self.__dict__['occupation'] = 'None'
38:      self.__dict__['profession'] = 'None'
39:      self.__dict__['industry'] = 'None'
40:      self.__dict__['sector'] = 'None'
41:      self.__dict__['field'] = 'None'
42:      self.__dict__['domain'] = 'None'
43:      self.__dict__['area'] = 'None'
44:      self.__dict__['specialty'] = 'None'
45:      self.__dict__['expertise'] = 'None'
46:      self.__dict__['skill'] = 'None'
47:      self.__dict__['talent'] = 'None'
48:      self.__dict__['ability'] = 'None'
49:      self.__dict__['aptitude'] = 'None'
50:      self.__dict__['capacity'] = 'None'
51:      self.__dict__['potential'] = 'None'
52:      self.__dict__['capability'] = 'None'
53:      self.__dict__['competence'] = 'None'
54:      self.__dict__['proficiency'] = 'None'
55:      self.__dict__['expertise'] = 'None'
56:      self.__dict__['mastery'] = 'None'
57:      self.__dict__['virtuosity'] = 'None'
58:      self.__dict__['genius'] = 'None'
59:      self.__dict__['brilliance'] = 'None'
60:      self.__dict__['intelligence'] = 'None'
61:      self.__dict__['wisdom'] = 'None'
62:      self.__dict__['knowledge'] = 'None'
63:      self.__dict__['understanding'] = 'None'
64:      self.__dict__['insight'] = 'None'
65:      self.__dict__['perception'] = 'None'
66:      self.__dict__['sight'] = 'None'
67:      self.__dict__['vision'] = 'None'
68:      self.__dict__['view'] = 'None'
69:      self.__dict__['opinion'] = 'None'
70:      self.__dict__['belief'] = 'None'
71:      self.__dict__['faith'] = 'None'
72:      self.__dict__['trust'] = 'None'
73:      self.__dict__['confidence'] = 'None'
74:      self.__dict__['assurance'] = 'None'
75:      self.__dict__['certainty'] = 'None'
76:      self.__dict__['conviction'] = 'None'
77:      self.__dict__['determination'] = 'None'
78:      self.__dict__['resolve'] = 'None'
79:      self.__dict__['will'] = 'None'
80:      self.__dict__['power'] = 'None'
81:      self.__dict__['strength'] = 'None'
82:      self.__dict__['force'] = 'None'
83:      self.__dict__['energy'] = 'None'
84:      self.__dict__['vitality'] = 'None'
85:      self.__dict__['vigor'] = 'None'
86:      self.__dict__['robustness'] = 'None'
87:      self.__dict__['hardiness'] = 'None'
88:      self.__dict__['toughness'] = 'None'
89:      self.__dict__['resilience'] = 'None'
90:      self.__dict__['endurance'] = 'None'
91:      self.__dict__['stamina'] = 'None'
92:      self.__dict__['constitution'] = 'None'
93:      self.__dict__['physique'] = 'None'
94:      self.__dict__['build'] = 'None'
95:      self.__dict__['frame'] = 'None'
96:      self.__dict__['structure'] = 'None'
97:      self.__dict__['form'] = 'None'
98:      self.__dict__['shape'] = 'None'
99:      self.__dict__['figure'] = 'None'
100:     self.__dict__['appearance'] = 'None'
101:     self.__dict__['look'] = 'None'
102:     self.__dict__['seemance'] = 'None'
103:     self.__dict__['physiognomy'] = 'None'
104:     self.__dict__['features'] = 'None'
105:     self.__dict__['traits'] = 'None'
106:     self.__dict__['characteristics'] = 'None'
107:     self.__dict__['qualities'] = 'None'
108:     self.__dict__['attributes'] = 'None'
109:     self.__dict__['properties'] = 'None'
110:     self.__dict__['characteristics'] = 'None'
111:     self.__dict__['features'] = 'None'
112:     self.__dict__['traits'] = 'None'
113:     self.__dict__['characteristics'] = 'None'
114:     self.__dict__['qualities'] = 'None'
115:     self.__dict__['attributes'] = 'None'
116:     self.__dict__['properties'] = 'None'
117:     self.__dict__['characteristics'] = 'None'
118:     self.__dict__['features'] = 'None'
119:     self.__dict__['traits'] = 'None'
120:     self.__dict__['characteristics'] = 'None'
121:     self.__dict__['qualities'] = 'None'
122:     self.__dict__['attributes'] = 'None'
123:     self.__dict__['properties'] = 'None'
124:     self.__dict__['characteristics'] = 'None'
125:     self.__dict__['features'] = 'None'
126:     self.__dict__['traits'] = 'None'
127:     self.__dict__['characteristics'] = 'None'
128:     self.__dict__['qualities'] = 'None'
129:     self.__dict__['attributes'] = 'None'
130:     self.__dict__['properties'] = 'None'
131:     self.__dict__['characteristics'] = 'None'
132:     self.__dict__['features'] = 'None'
133:     self.__dict__['traits'] = 'None'
134:     self.__dict__['characteristics'] = 'None'
135:     self.__dict__['qualities'] = 'None'
136:     self.__dict__['attributes'] = 'None'
137:     self.__dict__['properties'] = 'None'
138:     self.__dict__['characteristics'] = 'None'
139:     self.__dict__['features'] = 'None'
140:     self.__dict__['traits'] = 'None'
141:     self.__dict__['characteristics'] = 'None'
142:     self.__dict__['qualities'] = 'None'
143:     self.__dict__['attributes'] = 'None'
144:     self.__dict__['properties'] = 'None'
145:     self.__dict__['characteristics'] = 'None'
146:     self.__dict__['features'] = 'None'
147:     self.__dict__['traits'] = 'None'
148:     self.__dict__['characteristics'] = 'None'
149:     self.__dict__['qualities'] = 'None'
150:     self.__dict__['attributes'] = 'None'
151:     self.__dict__['properties'] = 'None'
152:     self.__dict__['characteristics'] = 'None'
153:     self.__dict__['features'] = 'None'
154:     self.__dict__['traits'] = 'None'
155:     self.__dict__['characteristics'] = 'None'
156:     self.__dict__['qualities'] = 'None'
157:     self.__dict__['attributes'] = 'None'
158:     self.__dict__['properties'] = 'None'
159:     self.__dict__['characteristics'] = 'None'
160:     self.__dict__['features'] = 'None'
161:     self.__dict__['traits'] = 'None'
162:     self.__dict__['characteristics'] = 'None'
163:     self.__dict__['qualities'] = 'None'
164:     self.__dict__['attributes'] = 'None'
165:     self.__dict__['properties'] = 'None'
166:     self.__dict__['characteristics'] = 'None'
167:     self.__dict__['features'] = 'None'
168:     self.__dict__['traits'] = 'None'
169:     self.__dict__['characteristics'] = 'None'
170:     self.__dict__['qualities'] = 'None'
171:     self.__dict__['attributes'] = 'None'
172:     self.__dict__['properties'] = 'None'
173:     self.__dict__['characteristics'] = 'None'
174:     self.__dict__['features'] = 'None'
175:     self.__dict__['traits'] = 'None'
176:     self.__dict__['characteristics'] = 'None'
177:     self.__dict__['qualities'] = 'None'
178:     self.__dict__['attributes'] = 'None'
179:     self.__dict__['properties'] = 'None'
180:     self.__dict__['characteristics'] = 'None'
181:     self.__dict__['features'] = 'None'
182:     self.__dict__['traits'] = 'None'
183:     self.__dict__['characteristics'] = 'None'
184:     self.__dict__['qualities'] = 'None'
185:     self.__dict__['attributes'] = 'None'
186:     self.__dict__['properties'] = 'None'
187:     self.__dict__['characteristics'] = 'None'
188:     self.__dict__['features'] = 'None'
189:     self.__dict__['traits'] = 'None'
190:     self.__dict__['characteristics'] = 'None'
191:     self.__dict__['qualities'] = 'None'
192:     self.__dict__['attributes'] = 'None'
193:     self.__dict__['properties'] = 'None'
194:     self.__dict__['characteristics'] = 'None'
195:     self.__dict__['features'] = 'None'
196:     self.__dict__['traits'] = 'None'
197:     self.__dict__['characteristics'] = 'None'
198:     self.__dict__['qualities'] = 'None'
199:     self.__dict__['attributes'] = 'None'
200:     self.__dict__['properties'] = 'None'
201:     self.__dict__['characteristics'] = 'None'
202:     self.__dict__['features'] = 'None'
203:     self.__dict__['traits'] = 'None'
204:     self.__dict__['characteristics'] = 'None'
205:     self.__dict__['qualities'] = 'None'
206:     self.__dict__['attributes'] = 'None'
207:     self.__dict__['properties'] = 'None'
208:     self.__dict__['characteristics'] = 'None'
209:     self.__dict__['features'] = 'None'
210:     self.__dict__['traits'] = 'None'
211:     self.__dict__['characteristics'] = 'None'
212:     self.__dict__['qualities'] = 'None'
213:     self.__dict__['attributes'] = 'None'
214:     self.__dict__['properties'] = 'None'
215:     self.__dict__['characteristics'] = 'None'
216:     self.__dict__['features'] = 'None'
217:     self.__dict__['traits'] = 'None'
218:     self.__dict__['characteristics'] = 'None'
219:     self.__dict__['qualities'] = 'None'
220:     self.__dict__['attributes'] = 'None'
221:     self.__dict__['properties'] = 'None'
222:     self.__dict__['characteristics'] = 'None'
223:     self.__dict__['features'] = 'None'
224:     self.__dict__['traits'] = 'None'
225:     self.__dict__['characteristics'] = 'None'
226:     self.__dict__['qualities'] = 'None'
227:     self.__dict__['attributes'] = 'None'
228:     self.__dict__['properties'] = 'None'
229:     self.__dict__['characteristics'] = 'None'
230:     self.__dict__['features'] = 'None'
231:     self.__dict__['traits'] = 'None'
232:     self.__dict__['characteristics'] = 'None'
233:     self.__dict__['qualities'] = 'None'
234:     self.__dict__['attributes'] = 'None'
235:     self.__dict__['properties'] = 'None'
236:     self.__dict__['characteristics'] = 'None'
237:     self.__dict__['features'] = 'None'
238:     self.__dict__['traits'] = 'None'
239:     self.__dict__['characteristics'] = 'None'
240:     self.__dict__['qualities'] = 'None'
241:     self.__dict__['attributes'] = 'None'
242:     self.__dict__['properties'] = 'None'
243:     self.__dict__['characteristics'] = 'None'
244:     self.__dict__['features'] = 'None'
245:     self.__dict__['traits'] = 'None'
246:     self.__dict__['characteristics'] = 'None'
247:     self.__dict__['qualities'] = 'None'
248:     self.__dict__['attributes'] = 'None'
249:     self.__dict__['properties'] = 'None'
250:     self.__dict__['characteristics'] = 'None'
251:     self.__dict__['features'] = 'None'
252:     self.__dict__['traits'] = 'None'
253:     self.__dict__['characteristics'] = 'None'
254:     self.__dict__['qualities'] = 'None'
255:     self.__dict__['attributes'] = 'None'
256:     self.__dict__['properties'] = 'None'
257:     self.__dict__['characteristics'] = 'None'

```

| | url | length_url | length_hostname | ip | nb_dots | nb_hyphens | nb_at | nb_qm | nb_and | nb_or | .. | domain_in_title | domains_with_copyright | whois_registered_domain | domain_registration_length |
|-------|-------|------------|-----------------|-------|---------|------------|-------|-------|--------|-------|----|-----------------|------------------------|-------------------------|----------------------------|
| 0 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 1 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 2 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 3 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 4 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| .. | .. | .. | .. | .. | .. | .. | .. | .. | .. | .. | .. | .. | .. | .. | .. |
| 11425 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 11426 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 11427 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 11428 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |
| 11429 | False | False | False | False | False | False | False | False | False | False | .. | False | False | False | False |

```
1]: missing_values=df.isnull().sum()
print("\n missing_values:")
missing_values>0
```

```
missing_values:
1]: url            False
length_url        False
length_hostname   False
ip                False
nb_dots           False
...
web_traffic       False
dns_record        False
google_index      False
page_rank         False
status           False
Length: 89, dtype: bool
```

```
1]: # Check for duplicate entries
duplicates = df.duplicated().sum()
print(f"\nNumber of duplicate rows: {duplicates}")
```

Number of duplicate rows: 0

```
: # Select numerical columns for outlier detection
numerical_cols = df.select_dtypes(include=['int64', 'float64']).columns.tolist()
numerical_cols.remove('page_rank')
```

```
: # Calculate IQR and remove outliers
for col in numerical_cols:
    Q1 = df[col].quantile(0.25)
    Q3 = df[col].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    df = df[(df[col] >= lower_bound) & (df[col] <= upper_bound)]

print("Outliers removed using IQR method.")
```

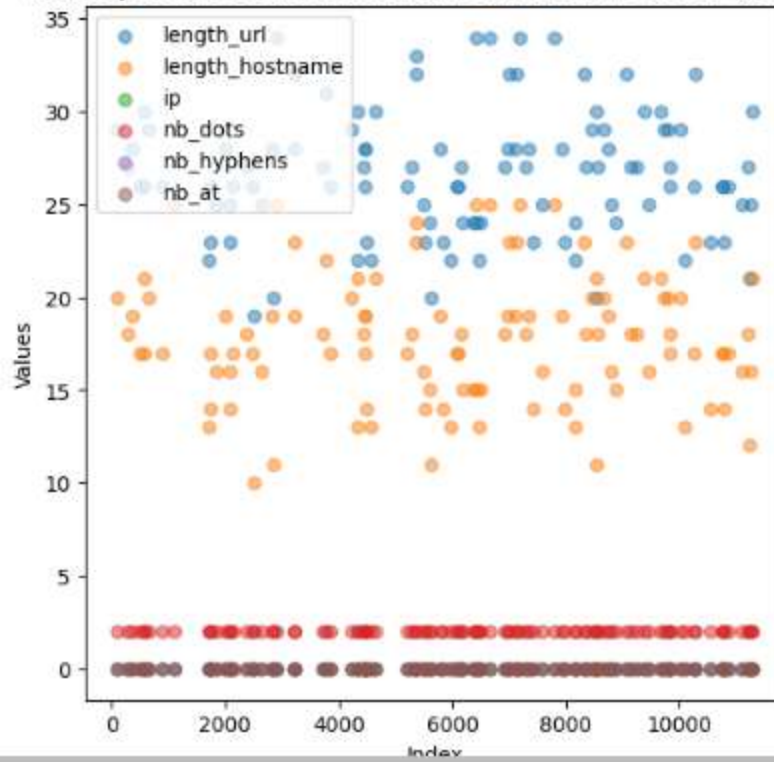
Outliers removed using IQR method.

```

1 # Scatterplot for outlier detection
plt.figure(figsize=(6, 6))
for col in numerical_cols[:6]:
    plt.scatter(df.index, df[col], label=col, alpha=0.5)
plt.legend()
plt.title("Scatterplot of Selected Numerical Features After Outlier Removal")
plt.ylabel("Values")
plt.xlabel("Index")
plt.show()

```

Scatterplot of Selected Numerical Features After Outlier Removal

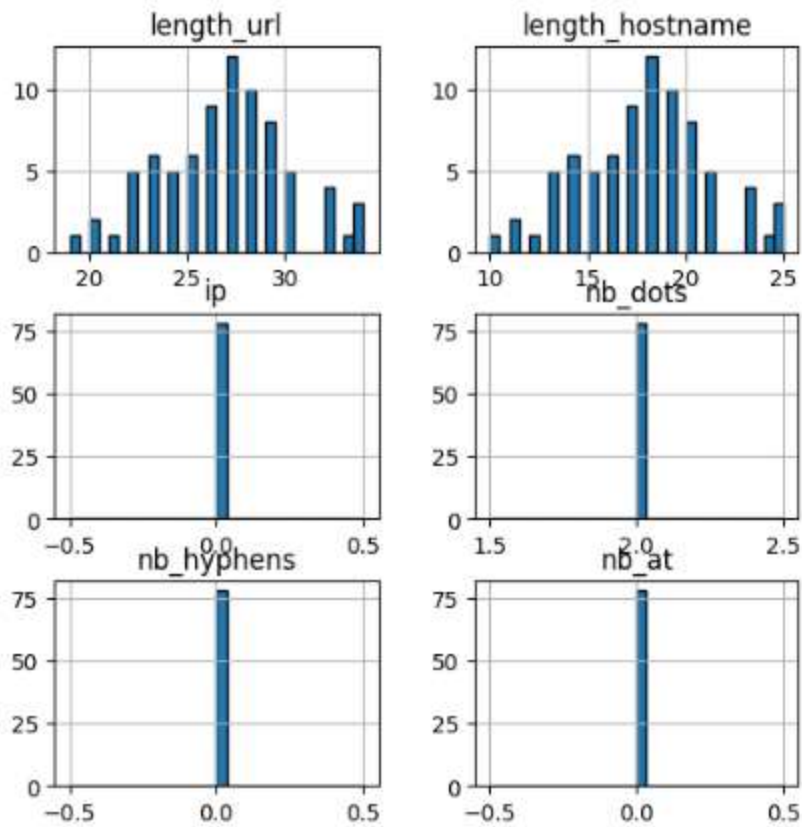


```

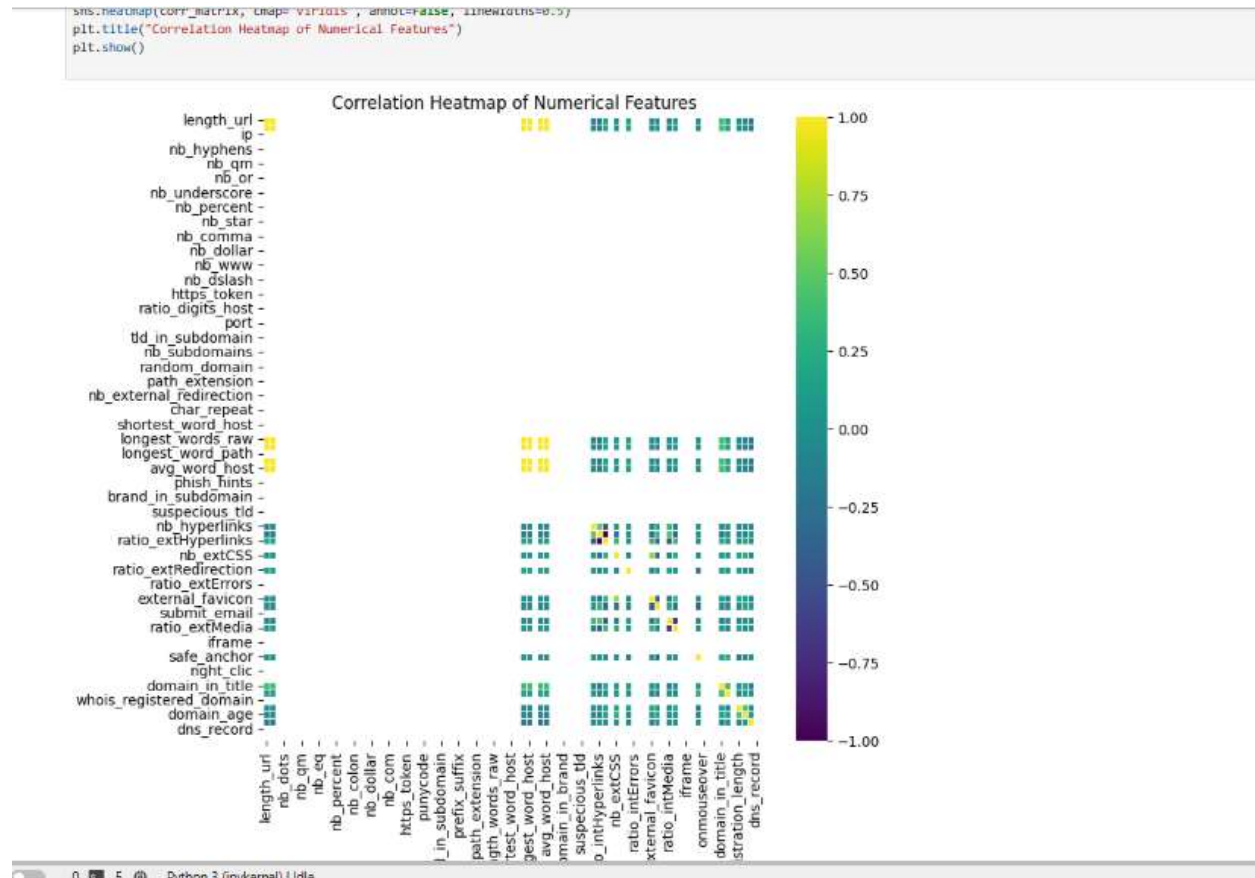
[: # Histograms for key numerical features
df[numerical_cols[:6]].hist(figsize=(6, 6), bins=30, edgecolor='black')
plt.suptitle("Histograms of Selected Numerical Features")
plt.show()

```

Histograms of Selected Numerical Features



Correlation heatmap:



Summary of the EDA Process in Theory

This Python script performs **Exploratory Data Analysis (EDA)** on a phishing website dataset using **Pandas, Matplotlib, and Seaborn**. The main steps involved are:

1. **Loading the Dataset**
 - Reads the dataset from a CSV file into a Pandas DataFrame.
 - Displays basic information such as column names, data types, and the first few rows.
2. **Handling Missing Values**
 - Checks for missing values in each column.
 - If missing values are found, it visualizes them using a bar plot.
 - If no missing values exist, it prints a confirmation message.

3. Outlier Detection and Removal Using IQR

- Identifies numerical columns (excluding `page_rank`).
- Uses the **Interquartile Range (IQR) method** to detect and remove outliers:
 - Computes the 1st quartile (Q1) and 3rd quartile (Q3).
 - Calculates the **IQR** as $Q3 - Q1$.
 - Defines the lower and upper bounds as $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$, respectively.
 - Filters out values beyond these bounds to remove outliers.

4. Visualization of Data Distributions

- **Boxplots**: Show the distribution of numerical features after outlier removal.
- **Histograms**: Display the frequency distribution of selected numerical features.

5. Correlation Analysis

- Computes the correlation matrix for numerical features.
- Plots a **heatmap** using the "**viridis**" colormap to visualize relationships between features.
- Helps identify highly correlated variables that may be redundant or important for classification.

Outcome of the EDA Process

- **Missing values**: Checked and handled if found.
- **Outliers**: Removed using IQR to improve data quality.
- **Data distribution**: Explored through histograms and boxplots.
- **Feature relationships**: Examined using a correlation heatmap.

This process ensures that the dataset is **cleaned, structured, and ready** for further analysis or machine learning modeling.