

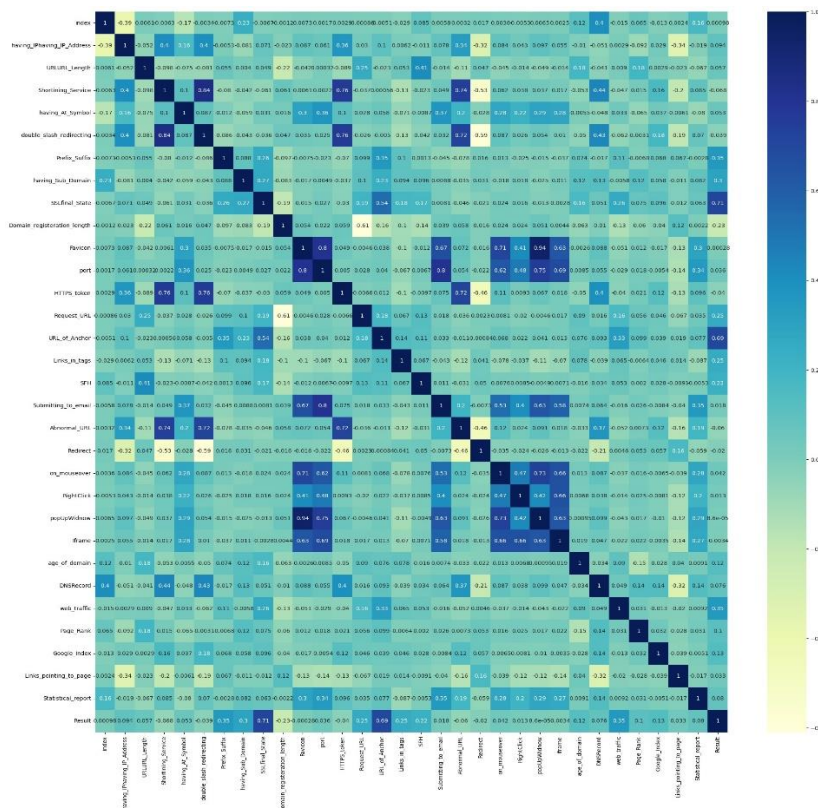
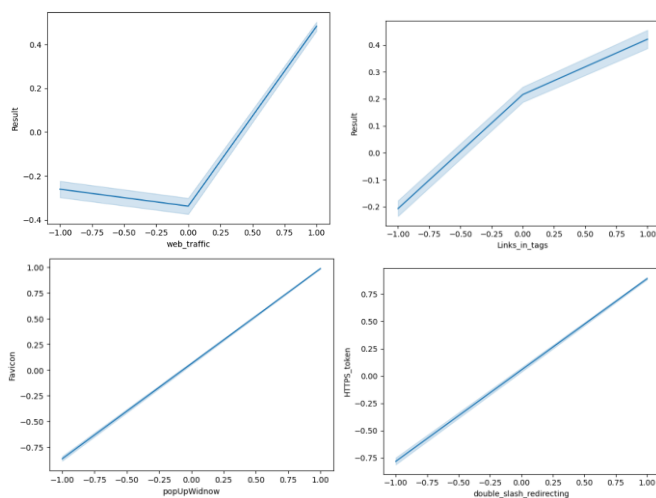
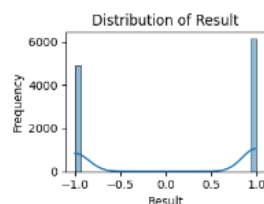
Data Collection and Preprocessing Phase

Date	29 September 2024
Team ID	LTVIP2024TMID24838
Project Title	Detection of Phishing Websites from URLs
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description																																																															
Data Overview	Dimensions: 11055 rows X 32 columns																																																															
	<table><thead><tr><th></th><th>index</th><th>having_IPhaving_IP_Address</th><th>URLURL_Length</th><th>Shortining_Service</th><th>having_At_Symbol</th><th>double_slash_redirecting</th></tr></thead><tbody><tr><td>count</td><td>11055.000000</td><td>11055.000000</td><td>11055.000000</td><td>11055.000000</td><td>11055.000000</td><td>11055.000000</td></tr><tr><td>mean</td><td>5528.000000</td><td>0.313795</td><td>-0.633198</td><td>0.738761</td><td>0.700588</td><td>0.741474</td></tr><tr><td>std</td><td>3191.447947</td><td>0.949534</td><td>0.766095</td><td>0.673998</td><td>0.713598</td><td>0.671011</td></tr><tr><td>min</td><td>1.000000</td><td>-1.000000</td><td>-1.000000</td><td>-1.000000</td><td>-1.000000</td><td>-1.000000</td></tr><tr><td>25%</td><td>2764.500000</td><td>-1.000000</td><td>-1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td></tr><tr><td>50%</td><td>5528.000000</td><td>1.000000</td><td>-1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td></tr><tr><td>75%</td><td>8291.500000</td><td>1.000000</td><td>-1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td></tr><tr><td>max</td><td>11055.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td><td>1.000000</td></tr></tbody></table>		index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474	std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011	min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000	25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000	50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000	75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000	max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000
		index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting																																																									
	count	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000	11055.000000																																																									
	mean	5528.000000	0.313795	-0.633198	0.738761	0.700588	0.741474																																																									
	std	3191.447947	0.949534	0.766095	0.673998	0.713598	0.671011																																																									
	min	1.000000	-1.000000	-1.000000	-1.000000	-1.000000	-1.000000																																																									
	25%	2764.500000	-1.000000	-1.000000	1.000000	1.000000	1.000000																																																									
	50%	5528.000000	1.000000	-1.000000	1.000000	1.000000	1.000000																																																									
75%	8291.500000	1.000000	-1.000000	1.000000	1.000000	1.000000																																																										
max	11055.000000	1.000000	1.000000	1.000000	1.000000	1.000000																																																										
Univariate Analysis	<div><div><p>Distribution of index</p></div><div><p>Distribution of having_IPhaving_IP_Address</p></div><div><p>Distribution of URLURL_Length</p></div><div><p>Distribution of Prefix_Suffix</p></div><div><p>Distribution of having_Sub_Domain</p></div><div><p>Distribution of SSLfinal_State</p></div></div>																																																															



Outliers and Anomalies	Identification and treatment of outliers.																																																																																																
Data Preprocessing Code Screenshots																																																																																																	
Loading Data	<div><pre>#Loading the Dataset df = pd.read_csv('/content/dataset_website (1).csv') df.head()</pre></div> <table><thead><tr><th></th><th>index</th><th>having_IPhaving_IP_Address</th><th>URLURL_Length</th><th>Shortining_Service</th><th>having_At_Symbol</th><th>double_slash_redirecting</th><th>Prefix_Suffix</th></tr></thead><tbody><tr><td>0</td><td>1</td><td>-1</td><td>1</td><td>1</td><td>1</td><td>-1</td><td>-1</td></tr><tr><td>1</td><td>2</td><td>1</td><td>1</td><td>1</td><td>1</td><td>1</td><td>-1</td></tr><tr><td>2</td><td>3</td><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td><td>-1</td></tr><tr><td>3</td><td>4</td><td>1</td><td>0</td><td>1</td><td>1</td><td>1</td><td>-1</td></tr><tr><td>4</td><td>5</td><td>1</td><td>0</td><td>-1</td><td>1</td><td>1</td><td>-1</td></tr></tbody></table>		index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix	0	1	-1	1	1	1	-1	-1	1	2	1	1	1	1	1	-1	2	3	1	0	1	1	1	-1	3	4	1	0	1	1	1	-1	4	5	1	0	-1	1	1	-1																																																
	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Prefix_Suffix																																																																																										
0	1	-1	1	1	1	-1	-1																																																																																										
1	2	1	1	1	1	1	-1																																																																																										
2	3	1	0	1	1	1	-1																																																																																										
3	4	1	0	1	1	1	-1																																																																																										
4	5	1	0	-1	1	1	-1																																																																																										
Handling Missing Data	<div><pre>df.isnull()</pre><div>✓ 0.0s</div><table><thead><tr><th></th><th>index</th><th>having_IPhaving_IP_Address</th><th>URLURL_Length</th><th>Shortining_Service</th><th>having_At_Symbol</th><th>double_slash_redirecting</th><th>Pref</th></tr></thead><tbody><tr><td>0</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>1</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>2</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>3</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>4</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>11050</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>11051</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>11052</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>11053</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr><tr><td>11054</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td><td>False</td></tr></tbody></table><div>11055 rows × 32 columns</div></div> <div>No NULL values found in the dataset.</div>		index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Pref	0	False	False	False	False	False	False	False	1	False	False	False	False	False	False	False	2	False	False	False	False	False	False	False	3	False	False	False	False	False	False	False	4	False	False	False	False	False	False	False	11050	False	False	False	False	False	False	False	11051	False	False	False	False	False	False	False	11052	False	False	False	False	False	False	False	11053	False	False	False	False	False	False	False	11054	False	False	False	False	False	False	False
	index	having_IPhaving_IP_Address	URLURL_Length	Shortining_Service	having_At_Symbol	double_slash_redirecting	Pref																																																																																										
0	False	False	False	False	False	False	False																																																																																										
1	False	False	False	False	False	False	False																																																																																										
2	False	False	False	False	False	False	False																																																																																										
3	False	False	False	False	False	False	False																																																																																										
4	False	False	False	False	False	False	False																																																																																										
...																																																																																										
11050	False	False	False	False	False	False	False																																																																																										
11051	False	False	False	False	False	False	False																																																																																										
11052	False	False	False	False	False	False	False																																																																																										
11053	False	False	False	False	False	False	False																																																																																										
11054	False	False	False	False	False	False	False																																																																																										
Data Transformation	No need for data transformation for the used dataset.																																																																																																
Feature Engineering	Included in the final code.																																																																																																
Save Processed Data	-																																																																																																