

# EMAIL PHISHING URL DETECTION

BEKZADAYEVA KYMBAT

SIS-2208

# INTRODUCTION

- PHISHING IS ONE OF THE MOST COMMON CYBERATTACKS, IN WHICH ATTACKERS IMPERSONATE LEGITIMATE ONLINE SERVICES IN ORDER TO DECEIVE USERS AND OBTAIN SENSITIVE INFORMATION. WITH THE RAPID GROWTH OF ONLINE PLATFORMS AND DIGITAL COMMUNICATION, PHISHING ATTACKS HAVE BECOME INCREASINGLY WIDESPREAD AND SOPHISTICATED.
- TRADITIONAL PHISHING DETECTION APPROACHES ARE MAINLY BASED ON BLACKLISTS AND MANUALLY DEFINED RULES. WHILE SUCH METHODS CAN SUCCESSFULLY BLOCK KNOWN THREATS, THEY ARE INEFFECTIVE AGAINST NEWLY GENERATED PHISHING URLS THAT CONSTANTLY CHANGE THEIR STRUCTURE TO EVADE DETECTION.
- MACHINE LEARNING PROVIDES AN EFFECTIVE SOLUTION TO THIS PROBLEM BY ENABLING SYSTEMS TO LEARN STRUCTURAL PATTERNS FROM DATA AND GENERALIZE TO PREVIOUSLY UNSEEN URLS. IN THIS PROJECT, A MACHINE LEARNING-BASED SYSTEM FOR PHISHING URL DETECTION IS DEVELOPED. THE MODEL CLASSIFIES URLS AS LEGITIMATE OR PHISHING USING URL-BASED FEATURES AND PROVIDES PROBABILITY-BASED PREDICTIONS. THE TRAINED MODEL IS DEPLOYED USING A FLASK WEB APPLICATION WITH A SIMPLE USER INTERFACE THAT ALLOWS REAL-TIME URL ANALYSIS AND ADJUSTABLE DECISION THRESHOLDS.

# DATASET AND EXPLORATORY ANALYSIS

The dataset used in this project was obtained from a public Kaggle repository called “Phishing Site URLs”. It originally contained over 549,000 URLs labeled as legitimate or phishing. After removing duplicate entries during preprocessing, approximately 507,000 unique URLs were used for training and evaluation. The dataset contains no missing values and is fully anonymized, making it suitable for academic research.

Since machine learning models require numerical target values, the original categorical labels were converted into binary format. Phishing URLs were encoded as 1, while legitimate URLs were encoded as 0, enabling supervised binary classification.

Exploratory data analysis showed that legitimate URLs dominate the dataset, while phishing URLs form a smaller but meaningful class, reflecting real-world conditions. Phishing URLs tend to be longer and contain more digits, special characters, and deeper path structures compared to legitimate URLs. These observations confirm that URL-level structural features provide strong signals for phishing detection and directly motivated the selected feature engineering approach.



← ↻ 🔍

https://www.kaggle.com/datasets/taruntiwarihp/phishing-site-urls

🏠 📁 ⭐ ⌂ 🔔 👤

Recently ... 🌐

☰ kaggle

+

Create

🏠

Home

🏆

Competitions

📁

Datasets

👤

Models

📊

Benchmarks

🎮

Game Arena

⌕

Code

💬

Discussions

🎓

Learn

⌵

More

📁

Your Work

⌵

viewed

📁

Phishing Site URLs

📁

Phishing URLs Datas...

📁

Web page Phishing D...

🔍

Search

👤

TARUN TIWARI · UPDATED 5 YEARS AGO

📈

174

⌕

Code

📄

Download

🔔

⋮

📄

Phishing Site URLs

Malicious and Phishing attacks uirs

📄

Data Card

⌕

Code (36)

💬

Discussion (0)

💡

Suggestions (0)

📄

About Dataset

Phishing site Predict dataset

[Youtube Explanation](#)

Content

- Data is containg 5,49,346 entries.
- There are two columns.
- Label column is prediction col which has 2 categories  
A. Good - which means the URLs is not containing malicious stuff and this site is not a Phishing Site.  
B. Bad - which means the URLs contains malicious stuff and this site is a Phishing Site.
- There is no missing value in the dataset.

Usability 📄

10.00

License

[Database: Open Database, Cont...](#)

Expected update frequency

Monthly

Tags

Internet

Social Networks

Crime

Public Safety

Search Engines

🖼️



# FEATURES.PY

As part of this project, a lexical feature extraction module was implemented in the file features.py. The purpose of this module is to transform a raw URL string into a fixed-length numerical feature vector suitable for machine learning classification.

The URL is first parsed using the urlparse function from the Python standard library to extract its main components, including the protocol, hostname, path, and query parameters. If a URL cannot be parsed correctly, a zero feature vector is returned to ensure system robustness and prevent runtime errors.

```
features.py X
features.py > ...
1  import re
2  import numpy as np
3  from urllib.parse import urlparse
4
5  # This function checks whether the hostname is an IP address
6  # Phishing URLs often use IP addresses instead of domain names
7  def has_ip_address(hostname: str) -> int:
8      pattern = r'^(\d{1,3}\.){3}\d{1,3}$'
9      return int(re.match(pattern, hostname) is not None)
10
11
12 # This function extracts numerical features from a URL string
13 # These features are later used by machine learning models
14 def extract_url_features(url: str) -> np.ndarray:
15     try:
16         # Parse the URL into its components (scheme, hostname, path, query, etc.)
17         parsed = urlparse(url)
18     except:
19         # If URL parsing fails, return a zero vector
20         return np.zeros(11)
21
22     # Extract hostname and path; use empty strings if missing
23     hostname = parsed.hostname or ""
24     path = parsed.path or ""
25
26     # Create a feature vector based on URL characteristics
27     return np.array([
28         len(url),                # Total length of the URL
29         len(hostname),           # Length of the hostname
30         len(path),               # Length of the path
31         sum(c.isdigit() for c in url), # Number of digits in the URL
32         sum(c in ['@', ':', '-', '?', '=', '%', '.', '#', '&'] for c in url), # Number of special characters
33         hostname.count('.'),      # Number of dots in the hostname
34         int('@' in url),          # Presence of '@' symbol
35         int('-' in hostname),     # Presence of '-' in hostname
36         int(parsed.scheme == "https"), # Whether HTTPS protocol is used
37         has_ip_address(hostname), # Whether hostname is an IP address
38         len(parsed.query)         # Length of query parameters
39     ], dtype=float)
```

# TF-IDF REPRESENTATION

- TF-IDF WITH CHARACTER-LEVEL N-GRAMS WAS USED TO CONVERT URLs INTO NUMERICAL FEATURE VECTORS SUITABLE FOR MACHINE LEARNING. CHARACTER-LEVEL REPRESENTATION IS PREFERRED OVER WORD-LEVEL TOKENIZATION BECAUSE URLs CONTAIN DIGITS, SYMBOLS, AND OBFUSCATED PATTERNS RATHER THAN NATURAL LANGUAGE.
- TF-IDF ASSIGNS HIGHER WEIGHTS TO DISTINCTIVE CHARACTER SEQUENCES THAT FREQUENTLY APPEAR IN PHISHING URLs WHILE REDUCING THE IMPACT OF COMMON PATTERNS. IN THIS PROJECT, N-GRAM LENGTHS FROM 3 TO 6 CHARACTERS WERE SELECTED TO BALANCE EXPRESSIVENESS AND COMPUTATIONAL EFFICIENCY. THIS APPROACH PROVIDES A ROBUST AND SCALABLE REPRESENTATION THAT GENERALIZES WELL TO PREVIOUSLY UNSEEN PHISHING URLs

# TRAIN\_MODEL.PY

```
train_model.py X
train_model.py > ...
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.svm import LinearSVC
5 from sklearn.calibration import CalibratedClassifierCV
6 from sklearn.metrics import (
7     classification_report,
8     confusion_matrix,
9     accuracy_score,
10    roc_auc_score,
11    average_precision_score
12 )
13 import joblib
14
15
16 # Path to the phishing URL dataset
17 DATA_PATH = r"C:\Users\Huawei\Documents\project\phishing_site_urls.csv\phishing_site_urls.csv"
18
19
20 # Load the dataset from CSV file into a DataFrame
21 df = pd.read_csv(DATA_PATH)
22
23
24 # Ensure consistent naming of the label column
25 # Some datasets use "Label", others use "label"
26 if "Label" in df.columns and "label" not in df.columns:
27     df = df.rename(columns={"Label": "label"})
28
29
30 # Remove rows with missing values and duplicate URLs
31 df = df.dropna(subset=["URL", "label"]).drop_duplicates(subset=["URL"])
32
33
34 # Convert all URLs to string type
35 df["URL"] = df["URL"].astype(str)
36
```

In the `train_model.py` module, I implemented the complete training and evaluation pipeline for phishing URL classification. The dataset is loaded from a CSV file using the pandas library and contains URLs with corresponding phishing or legitimate labels.

During preprocessing, I ensured data consistency by standardizing the label column name, removing rows with missing values, and eliminating duplicate URLs to avoid biased evaluation. All URLs were explicitly converted to string format to ensure compatibility with text-based feature extraction methods.

These preprocessing steps improve data quality and ensure that the model is trained on clean, consistent, and reliable input data before feature extraction and model training.



```

37
38 # Convert textual labels into binary numerical labels
39 # 1 represents phishing, 0 represents legitimate
40 def encode_label(x):
41     x = str(x).strip().lower()
42     return 1 if x in ["bad", "phishing", "malicious", "defacement", "1"] else 0
43
44
45 # Apply label encoding to the dataset
46 df["label"] = df["label"].apply(encode_label).astype(int)
47
48
49 # Split the dataset into training and testing sets
50 # Stratification preserves the class distribution
51 X_train, X_test, y_train, y_test = train_test_split(
52     df["URL"],
53     df["label"].values,
54     test_size=0.2,
55     random_state=42,
56     stratify=df["label"].values
57 )
58
59
60 # Initialize TF-IDF vectorizer using character-level n-grams
61 # Character n-grams are effective for detecting phishing patterns in URLs
62 vectorizer = TfidfVectorizer(
63     analyzer="char",
64     ngram_range=(3, 6),
65     min_df=2,
66     max_features=300000
67 )
68
69
70 # Fit the vectorizer on training data and transform both sets
71 X_train_vec = vectorizer.fit_transform(X_train)
72 X_test_vec = vectorizer.transform(X_test)
73

```

- In this stage, I converted the original textual labels into binary numerical values, where phishing URLs are encoded as 1 and legitimate URLs as 0. This enables supervised binary classification and ensures compatibility with standard machine learning models and evaluation metrics.
- To evaluate the model's generalization ability, the dataset was split into training and testing sets using an 80/20 ratio with stratification. Stratified splitting preserves the original class distribution in both sets, while a fixed random seed ensures reproducibility and fair evaluation without data leakage.
- After the split, URLs were transformed into numerical feature vectors using character-level TF-IDF n-grams with lengths from 3 to 6. The TF-IDF vectorizer was fitted only on the training data and then applied to both training and test sets, ensuring that no information from the test set was used during training.



```

train_model.py X
train_model.py > ...

75 # Create a Linear Support Vector Machine classifier
76 # Class balancing compensates for class imbalance in phishing datasets
77 base = LinearSVC(class_weight="balanced")
78
79
80 # Calibrate the SVM classifier to enable probability estimation
81 model = CalibratedClassifierCV(
82     base,
83     method="sigmoid",
84     cv=3
85 )
86
87
88 # Train the calibrated model on the training data
89 model.fit(X_train_vec, y_train)
90
91
92 # Predict phishing probabilities for the test set
93 proba = model.predict_proba(X_test_vec)[: , 1]
94
95
96 # Convert probabilities into binary predictions using threshold 0.5
97 pred = (proba >= 0.5).astype(int)
98
99
100 # Print evaluation metrics
101 print("Accuracy:", accuracy_score(y_test, pred))
102 print("ROC-AUC:", roc_auc_score(y_test, proba))
103 print("PR-AUC:", average_precision_score(y_test, proba))
104 print(classification_report(y_test, pred, digits=4))
105 print(confusion_matrix(y_test, pred))
106
107
108 # Save the trained model and TF-IDF vectorizer to disk
109 # These files are later used in the web application
110 joblib.dump(model, "model.pkl")
111 joblib.dump(vectorizer, "vectorizer.pkl")
112
113
114 print("Saved: model.pkl")
115 print("Saved: vectorizer.pkl")
116

```

- For classification, a Linear Support Vector Machine was selected because it performs well on high-dimensional sparse TF-IDF features. Class weighting was applied to address imbalance between phishing and legitimate URLs.
- Since LinearSVC does not provide probability estimates, the model was calibrated using CalibratedClassifierCV to produce phishing probabilities for threshold-based decision making.
- The model was evaluated on a held-out test set using accuracy, ROC-AUC, PR-AUC, and the confusion matrix, demonstrating strong performance without data leakage.

# MODEL EVALUATION RESULTS

After training, the model was evaluated on a held-out test set. The achieved accuracy is **98.15%**, indicating a high overall correctness of predictions. The ROC-AUC score of 0.9966 demonstrates excellent discriminative ability between phishing and legitimate URLs, while the PR-AUC score of **0.9913** confirms strong performance under class imbalance.

The classification report shows balanced precision and recall for both classes, and the confusion matrix indicates a low number of misclassifications. Importantly, these results were obtained without data leakage, as feature extraction and model training were performed only on the training set, and evaluation was conducted on unseen test data.

```
PROBLEMS  OUTPUT  DEBUG CONSOLE  TERMINAL  PORTS

PS C:\Users\Huawei\Documents\project> python train_model.py
Accuracy: 0.9814568361281164
ROC-AUC: 0.9966156961230328
PR-AUC: 0.991298218186387

              precision    recall  f1-score   support

         0       0.9846      0.9915      0.9881       78579
         1       0.9701      0.9469      0.9584       22860

   accuracy                   0.9815       101439
  macro avg              0.9774      0.9692      0.9732       101439
weighted avg              0.9814      0.9815      0.9814       101439

[[77913   666]
 [ 1215 21645]]
Saved: model.pkl
Saved: vectorizer.pkl
PS C:\Users\Huawei\Documents\project> 
```

# APP.PY

- After training the model, I deployed it as a web application using the Flask framework. The application loads the trained machine learning model and the TF-IDF vectorizer from disk using joblib and performs real-time phishing URL classification.
- Users can submit a URL through the web interface, after which the system predicts the phishing probability and converts it into a human-readable risk level (Low, Medium, or High). Loading the model once at application startup ensures efficient inference and avoids retraining for each request.

```
app.py X
app.py > ...
1  from flask import Flask, render_template, request
2  import joblib
3  # Create Flask application instance
4  app = Flask(__name__)
5  # Load the trained machine learning model from disk
6  # This model was trained earlier using TF-IDF features
7  model = joblib.load("model.pkl")
8  # Load the TF-IDF vectorizer used during training
9  # It converts a URL string into numerical feature vectors
10 vectorizer = joblib.load("vectorizer.pkl")
11
12 def risk_level(p):
13     #Convert phishing probability into a human-readable risk level.
14     # This improves interpretability for end users.
15     if p >= 0.75:
16         return "High"
17     if p >= 0.40:
18         return "Medium"
19     return "Low"
20
21 @app.route("/", methods=["GET", "POST"])
22 def index():
23     #Main route of the web application.
24     #Handles both displaying the form (GET)
25     #and processing user input (POST).
26
27     # Initialize variables passed to the HTML template
28     result = None
29     probability = None
30     risk = None
31     url_input = ""
32     threshold = 0.5 # Default decision threshold
33
```

```
34
35 # If the user submits the form
36 if request.method == "POST":
37
38     # Read URL entered by the user
39     url_input = request.form.get("url", "").strip()
40     # Read threshold value entered by the user
41     t = request.form.get("threshold", "0.5").strip()
42     # Try to convert threshold to float
43     try:
44         threshold = float(t)
45     except:
46         threshold = 0.5
47
48     # If a URL was provided
49     if url_input:
50         # Convert the URL into TF-IDF feature vector
51         x = vectorizer.transform([url_input])
52         # Predict phishing probability using the trained model
53         proba = float(model.predict_proba(x)[0][1])
54         # Apply decision threshold to get final class
55         pred = 1 if proba >= threshold else 0
56         # Convert numerical prediction to label
57         result = "Phishing" if pred == 1 else "Legitimate"
58
59         # Save probability and risk level for display
60         probability = proba
61         risk = risk_level(proba)
62
63     # Render HTML template with prediction results
64     return render_template("index.html", result=result, probability=probability, risk=risk, url_input=url_input, threshold=threshold)
65
66 # Run Flask development server
67 if __name__ == "__main__":
68     app.run(debug=True)
```



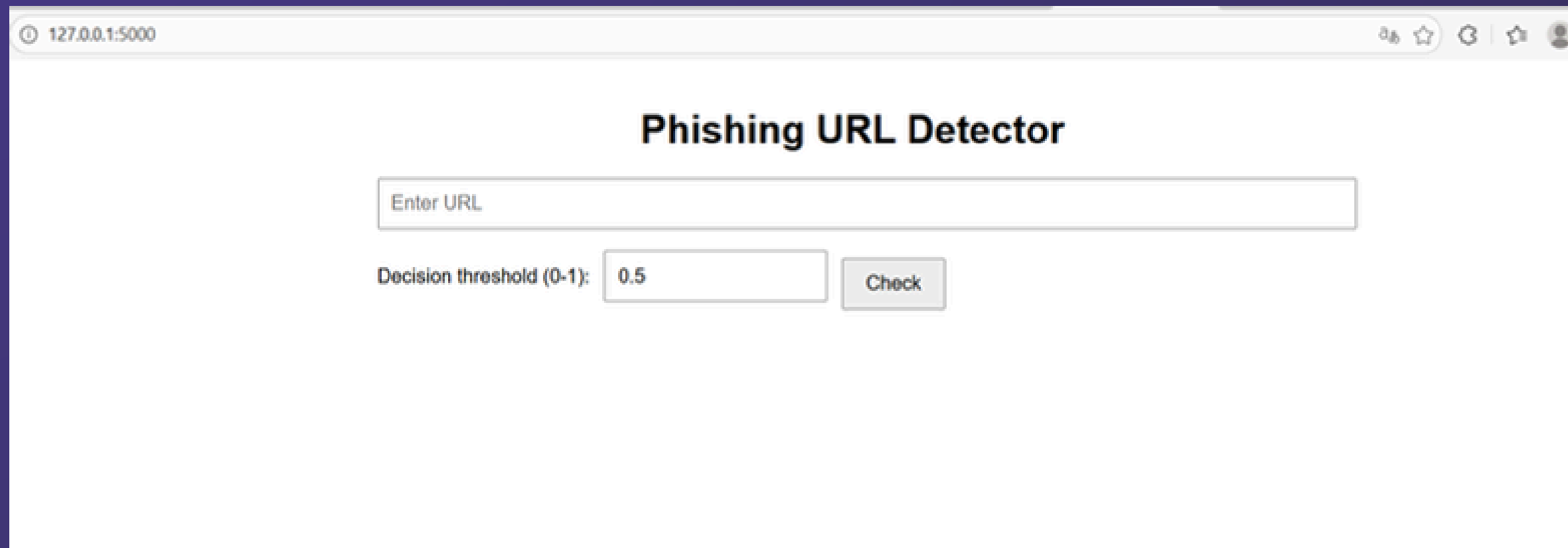
# USER INTERFACE

- After model inference, the application renders the index.html template and displays the prediction results to the user. The frontend receives the predicted class label, phishing probability, risk level, submitted URL, and the selected decision threshold from the backend.
- The user interface is implemented using HTML with embedded CSS and follows a simple, minimalistic design. It allows users to submit a URL, adjust the decision threshold, and view results only after a prediction is made, ensuring clarity and ease of use. This separation between backend logic and frontend presentation makes the application clean, interpretable, and easy to maintain.

```
templates > index.html > ...
1  <!DOCTYPE html>
2  <html lang="en">
3  <head>
4    <meta charset="UTF-8">
5    <title>Phishing URL Detector</title>
6    <style>
7      body { font-family: Arial, sans-serif; max-width: 780px; margin: 40px auto; }
8      input { width: 100%; padding: 10px; font-size: 16px; }
9      button { margin-top: 12px; padding: 10px 18px; font-size: 16px; cursor: pointer; }
10     .row { display: flex; gap: 12px; margin-top: 12px; align-items: center; }
11     .row input { width: 160px; }
12     .box { margin-top: 18px; padding: 14px; border-radius: 10px; }
13     .good { background: #e9f7ef; border: 1px solid #b7e1cd; }
14     .bad { background: #fdecea; border: 1px solid #f5c6cb; }
15   </style>
16 </head>
17 <body>
18   <h1 style="text-align:center;">Phishing URL Detector</h1>
19
20   <form method="POST">
21     <input type="text" name="url" placeholder="Enter URL" value="{{ url_input }}">
22     <div class="row">
23       <div>Decision threshold (0-1):</div>
24       <input type="text" name="threshold" value="{{ threshold }}">
25       <button type="submit">Check</button>
26     </div>
27   </form>
28
29   {% if result %}
30   <div class="box {% if result == 'Legitimate' %}good{% else %}bad{% endif %}">
31     <div><b>Result:</b> {{ result }}</div>
32     <div><b>Phishing probability:</b> {{ (probability * 100) | round(2) }}%</div>
33     <div><b>Risk level:</b> {{ risk }}</div>
34   </div>
35   {% endif %}
36 </body>
37 </html>
38
```

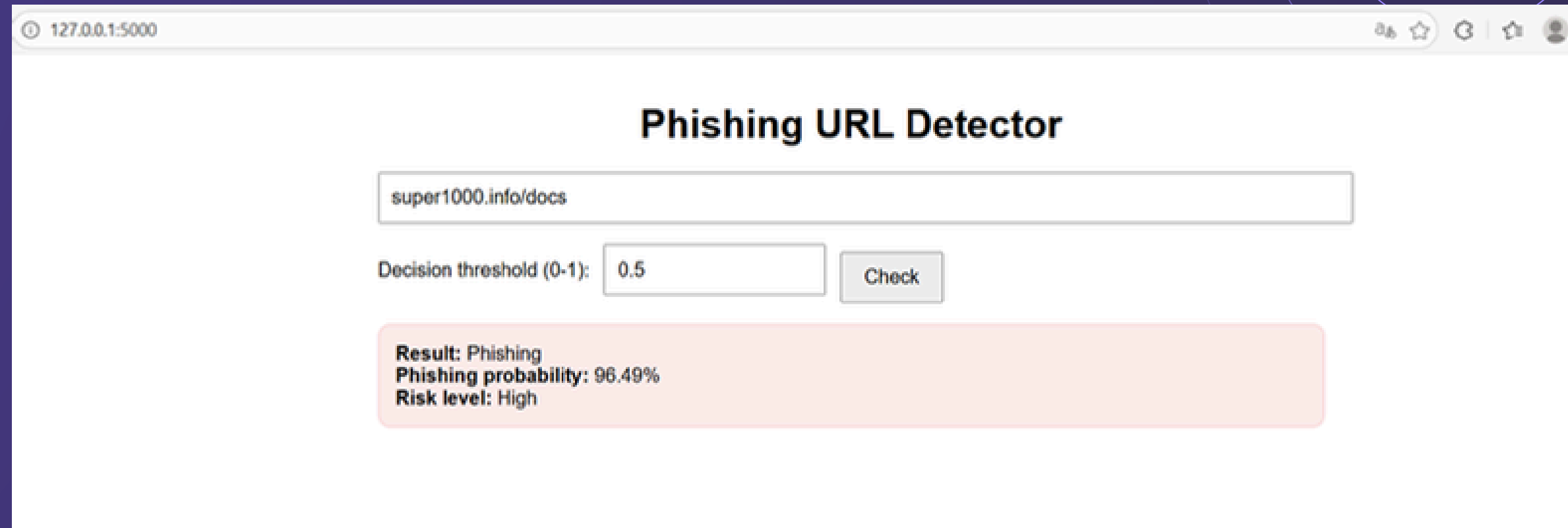
# USER INTERFACE OVERVIEW

The web application provides a simple interface for phishing URL detection. Users can enter a URL for analysis and adjust the decision threshold to control detection sensitivity. The interface follows a minimalistic design to ensure clarity and ease of use. It directly connects user input with the deployed machine learning model for real-time predictions.



The screenshot shows a web browser window with the address bar displaying "127.0.0.1:5000". The page title is "Phishing URL Detector". The interface is minimalist and white. It features a large text input field labeled "Enter URL". Below this, there is a label "Decision threshold (0-1):" followed by a smaller input field containing the value "0.5". To the right of the threshold input is a button labeled "Check".

The figure shows an example of the system's output when analyzing a suspicious URL. The model classifies the URL as phishing with a predicted probability of 96.49%, corresponding to a high risk level.

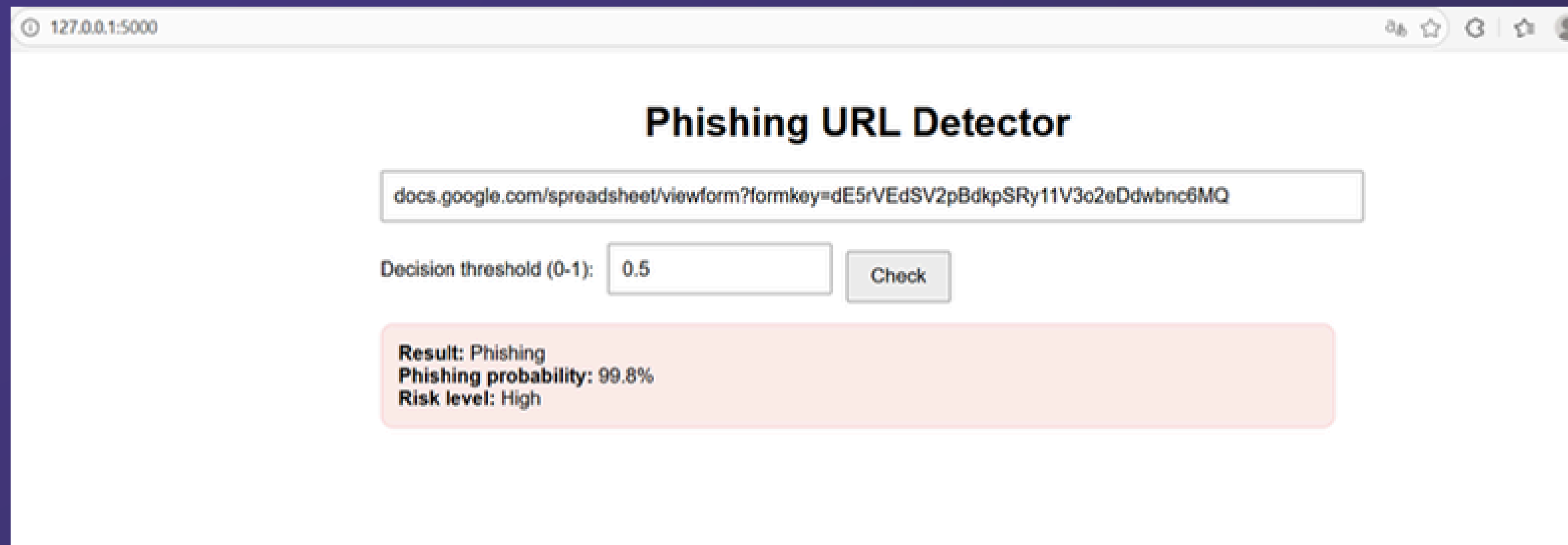


The screenshot displays a web browser window with the address bar showing '127.0.0.1:5000'. The page title is 'Phishing URL Detector'. Below the title, there is a text input field containing the URL 'super1000.info/docs'. Underneath the input field, the text 'Decision threshold (0-1):' is followed by a text input field containing '0.5' and a 'Check' button. Below these elements, a red-bordered box contains the following text: 'Result: Phishing', 'Phishing probability: 96.49%', and 'Risk level: High'.

| Field                    | Value               |
|--------------------------|---------------------|
| URL                      | super1000.info/docs |
| Decision threshold (0-1) | 0.5                 |
| Result                   | Phishing            |
| Phishing probability     | 96.49%              |
| Risk level               | High                |

The result is clearly presented using the predicted label, probability score, and risk interpretation, with a red background highlighting the threat. This example demonstrates successful real-time integration of the trained machine learning model into the web application.

The next figure illustrates another example of the system's output when analyzing a suspicious URL. The model classifies the URL as phishing with a predicted probability of 99.8%, corresponding to a high risk level. The result demonstrates the model's ability to confidently identify highly suspicious URLs using probability-based classification. Visual highlighting emphasizes the severity of the detected threat and improves interpretability of the output.



The screenshot shows a web browser window with the address bar displaying '127.0.0.1:5000'. The page title is 'Phishing URL Detector'. Below the title, there is a text input field containing the URL 'docs.google.com/spreadsheet/viewform?formkey=dE5rVEdSV2pBdkpSRy11V3o2eDdwbn6MQ'. Below the input field, there is a label 'Decision threshold (0-1):' followed by a text input field containing '0.5' and a 'Check' button. Below the 'Check' button, there is a red-bordered box containing the following text: 'Result: Phishing', 'Phishing probability: 99.8%', and 'Risk level: High'.

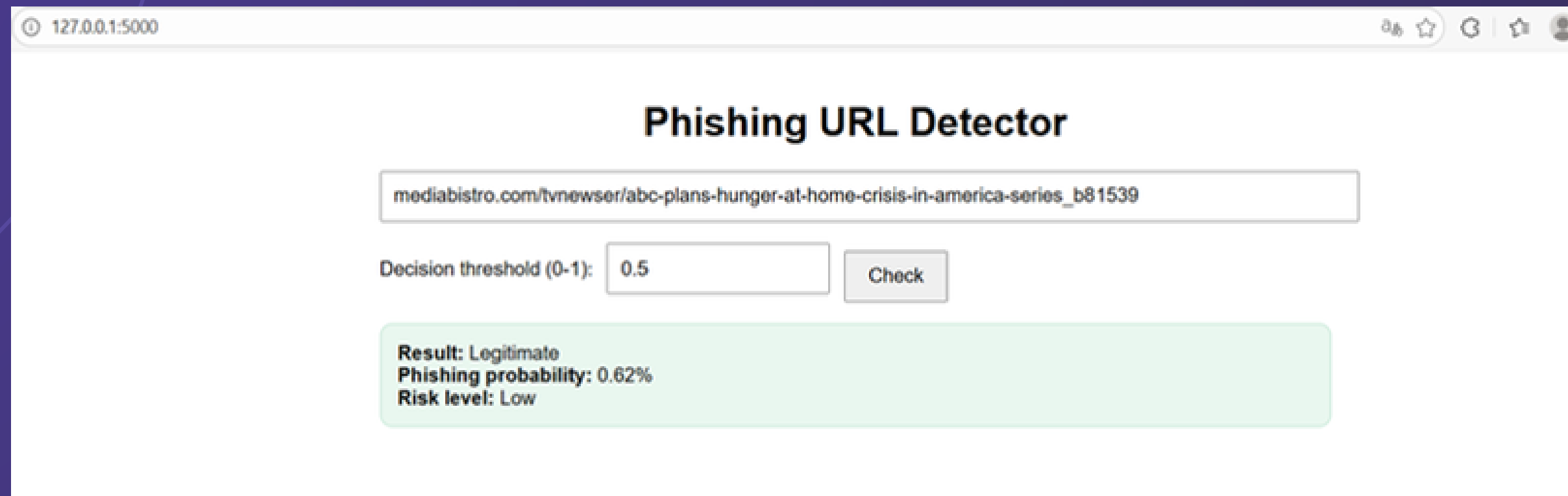
**Phishing URL Detector**

docs.google.com/spreadsheet/viewform?formkey=dE5rVEdSV2pBdkpSRy11V3o2eDdwbn6MQ

Decision threshold (0-1): 0.5

**Result: Phishing**  
**Phishing probability: 99.8%**  
**Risk level: High**





The figure presents an example of the system's output when analyzing a legitimate URL. The model classifies the URL as legitimate with a phishing probability of 0.62%, corresponding to a low risk level. This result demonstrates the model's ability to correctly identify benign URLs and avoid false positives. Together with the phishing detection examples, this confirms that the system provides balanced and reliable real-time predictions for both classes.

# CONCLUSION

- IN THIS PROJECT, A MACHINE LEARNING-BASED SYSTEM FOR PHISHING URL DETECTION WAS DEVELOPED USING CHARACTER-LEVEL TF-IDF FEATURES AND A CALIBRATED LINEAR SUPPORT VECTOR MACHINE CLASSIFIER. THE PROPOSED APPROACH SUCCESSFULLY DISTINGUISHES BETWEEN PHISHING AND LEGITIMATE URLS USING ONLY URL-LEVEL INFORMATION.
- A COMPLETE MACHINE LEARNING PIPELINE WAS IMPLEMENTED, INCLUDING DATA PREPROCESSING, FEATURE ENGINEERING, MODEL TRAINING, EVALUATION, AND DEPLOYMENT. CHARACTER-LEVEL TF-IDF N-GRAMS PROVED EFFECTIVE FOR CAPTURING STRUCTURAL PATTERNS ASSOCIATED WITH PHISHING ATTACKS AND ACHIEVED STRONG PERFORMANCE WITHOUT DATA LEAKAGE.
- THE TRAINED MODEL WAS DEPLOYED AS A FLASK WEB APPLICATION THAT PROVIDES REAL-TIME PHISHING DETECTION WITH PROBABILITY-BASED RISK ASSESSMENT. THE RESULTS CONFIRM THE EFFECTIVENESS OF MACHINE LEARNING FOR URL-BASED PHISHING DETECTION, WITH FUTURE WORK FOCUSING ON ADDITIONAL FEATURES OR ADVANCED MODELS.

**THANK YOU FOR**  
*attention*