

**DR. VISHWANATH KARAD MIT WORLD PEACE
UNIVERSITY, PUNE**

**Artificial Intelligence and Machine Learning Techniques
Third Year B. Tech, Semester 5**

**NEED FOR SECURITY - MOST WANTED
OSINT WITH AI AND ML**

MINI PROJECT REPORT

**Under the Guidance of
Dr. Yogita Hande**

Prepared By

Krishnaraj Thadesar, PA10, 1032210888
Sourab Karad, PA25, 1032211150
Saubhagya Singh, PA24, 1032211144
Parth Zarekar, PA06, 1032210846

**Department of School of Computer Engineering and
Technology
Maharashtra, India.
2023-2024
December 12, 2023**

Contents

1	Introduction	1
1.0.1	Problem Statement	1
1.0.2	Need of the Project	1
2	Literature Survey	2
2.1	Tools researched	2
2.1.1	CFLW	2
2.1.2	pimEyes	2
2.1.3	Predicta search	3
2.1.4	OSINT Industries	3
2.1.5	BreachDirectory	4
3	Methodology, Algorithms and Implementations	6
3.1	Methodology	6
3.2	Setup	6
3.3	Data Collection	6
3.3.1	Google Dorking	6
3.3.2	Web Scraping with Python	6
3.4	Feature Extraction	7
3.5	Analysis	7
3.6	Limitations	8
3.6.1	Ethical and Legal Compliance	8
3.7	Considerations	8
3.8	Algorithms	8
3.8.1	Naïve Bayes	8
3.8.2	K-Means	9
3.8.3	YACA	10
3.9	Implementation	11
3.9.1	Data For Training	11
3.9.2	Web Scrapping and Dorking Results	11
3.9.3	Training Results	13
3.10	Platform	13
3.11	Screenshots	14
4	Conclusion, Result and Discussions	17
5	Future Prospects	18
Bibliography		19

Acknowledgment

I would like to express my deepest appreciation to all those who provided me the possibility to complete this report. A special gratitude I give to our mentor, Dr. Yogita Hande, whose contribution in stimulating suggestions and encouragement, helped me to coordinate my project especially in writing this report.

Furthermore, I would also like to acknowledge with much appreciation the crucial role of the staff of MIT WPU, who gave the permission to use all required equipment and the necessary materials to complete the task. A special thanks goes to my team mates, who helped me enormously to assemble the parts and gave suggestion about the task of using the techniques of measurements.

I have to appreciate the guidance given by other supervisor as well as the panels especially in our project presentation that has improved our presentation skills thanks to their comment and advices.

I would also like to thank my parents for their wise counsel and sympathetic ear. You are always there for me. Finally, I wish to thank my friends for their support and encouragement throughout my study.

Name of Students

1. Krishnaraj Thadesar, PA10, 1032210888
2. Sourab Karad, PA25, 1032211150
3. Saubhagya Singh, PA24, 1032211144
4. Parth Zarekar, PA06, 1032210846

Abstract

In an increasingly digital world, the impact of our online presence is often overlooked. This project introduces NFS-most wanted, a platform that analyzes and evaluates our digital footprint. NFS-most wanted provides clear insights derived from our online interactions, simplifying the complexities and offering easily understandable observations.

Think of NFS-most wanted as a virtual mirror that reflects our digital identity. The user-friendly experience eliminates technical jargon and presents meaningful insights. Whether we want to understand how others perceive us online or take control of our digital identity, NFS-most wanted aims to provide a snapshot of our online persona.

The problem statement of this project is to build a website that analyzes and evaluates our online presence in today's interconnected world. Our digital footprint plays a significant role in shaping how others perceive us, yet many individuals are unaware of its extent and the potential risks associated with it.

The objective of this project is to develop a comprehensive platform that provides insights and awareness about our online presence. By analyzing and evaluating the information available about us on the internet, the website aims to help individuals understand the impact of their digital footprint and take necessary measures to manage and protect their online identity.

This project addresses the need for individuals to have a better understanding of their online presence and the importance of managing their digital identity. NFS-most wanted aims to empower users to make informed decisions about their online activities and take control of their digital footprint.

List of Figures

2.1	CFLW	2
2.2	pimEyes	3
2.3	predicta search	3
2.4	OSINT Industries	4
2.5	OSINT Industries Pricing	4
2.6	BreachDirectory	5
3.1	Example of Google Dorking	7
3.2	CSV File for Training Profession Prediction	11
3.3	Training Results from Backend	13
3.4	Home page	14
3.5	Form	14
3.6	Results	15
3.7	Key Features	15
3.8	Vision	16
3.9	Team	16

Chapter 1

Introduction

In an increasingly digital world, have you ever considered the impact of your online presence? Introducing NFS-most wanted, a comprehensive platform for a quick and straightforward analysis of your digital footprint. Here, we simplify the complexities and provide clear insights derived from your online interactions.

Think of NFS-most wanted as a virtual mirror that reflects your digital identity. We have designed a user-friendly experience that eliminates technical jargon, offering you easily understandable observations extracted from your online activities.

Our mission is to provide you with a snapshot of your online persona in the vast landscape of the internet. We focus on presenting meaningful insights, whether you are interested in understanding how others perceive you online or taking control of your digital identity.

1.0.1 Problem Statement

The problem statement of the project is to build a website that analyzes and evaluates our online presence in today's interconnected world. Our digital footprint, whether for personal or professional reasons, plays a significant role in shaping how others perceive us. However, many individuals are unaware of the extent of their digital footprint and the potential risks associated with it.

The objective of this project is to develop a comprehensive platform that provides insights and awareness about our online presence. By analyzing and evaluating the information available about us on the internet, the website aims to help individuals understand the impact of their digital footprint and take necessary measures to manage and protect their online identity.

1.0.2 Need of the Project

In today's interconnected world, our online presence plays a significant role in shaping how we are perceived by others. Whether for personal or professional reasons, it has become crucial to understand and manage the information available about us on the internet. However, many individuals are unaware of the extent of their digital footprint and the potential risks associated with it.

1. **Digital Footprint Unveiled:** Individuals often underestimate the information available about them online, scattered across various platforms. Lack of awareness regarding the implications of their digital footprint can lead to privacy concerns and potential security threats.
2. **Online Reputation Management:** Building and maintaining a positive online presence is vital for personal and professional success. Without proper tools and insights, individuals may struggle to curate their online image effectively.
3. **Privacy and Security Risks:** In an era of increasing cyber threats, understanding potential vulnerabilities is crucial. Lack of awareness regarding one's online vulnerabilities can result in identity theft, scams, or other cybercrimes.

Chapter 2

Literature Survey

In the literature survey, we reviewed various tools related to our project or to train our model which works in that specific manner. We explored different research papers, articles, and case studies to gain insight into this state of the art in this field. the systems provided valuable information on the techniques used other sources.

2.1 Tools researched

2.1.1 CFLW

CFLW Cyber Strategies (CFLW) is a company that provides cyber security services. CFLW's services include: Strategic studies, Capacity plans, Dialogues, Data analytics, Artificial intelligence.

CFLW Intelligence Services are based on in-depth expertise in Dark Web, Crypto-assets, Blockchain, Distributed Cryptography and Artificial Intelligence.



Intelligence Services
Strategic Insights, Operational Perspectives

CFLW Cyber Strategies provides solutions at the intersection of strategy and technology. Drawing on our well-established technical foundation and global presence, we provide intelligence services based on these solutions:

Figure 2.1: CFLW

2.1.2 pimEyes

pimEyes is a company that provides facial recognition services. pimEyes's services include: Facial recognition, Facial recognition, Facial recognition, Facial recognition, Facial recognition.

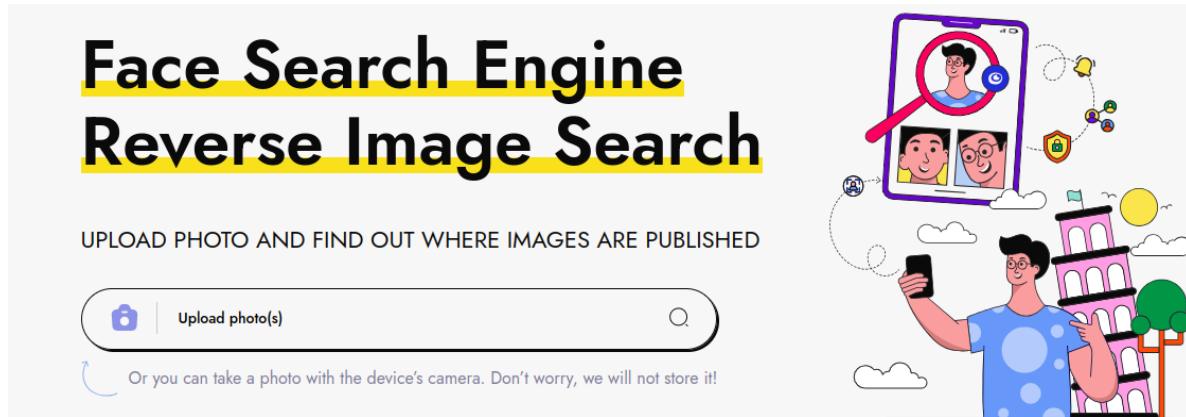


Figure 2.2: pimEyes

2.1.3 Predicta search

Predicta Search is a tool that provides related social media profiles when given an email address or phone number. It can be used for investigations or to get a digital footprint.

The screenshot shows the search results for the email address "parthzarekar@gmail.com". At the top, there are three buttons: "Show all 19", "Show found 4", and "Show not found 11". On the right, it says "Search is completed" with a green checkmark and "19 / 19". Below this, a message says "Search is completed and a detailed report has been generated" with a "Access report" button. The main area contains 19 search results in a grid format, each with a checked checkbox and a platform name. Most platforms show "No account found". Some platforms like GitHub, Notion, and Trello have "SIGN UP" buttons. The platforms listed are: About.me, Askfm, Flickr, Garmin, Gravatar, ImageShack, Notion, Picsart, Runkeeper, Skype, and Trello.

Figure 2.3: predicta search

2.1.4 OSINT Industries

OSINT Industries is a platform that uses open-source intelligence (OSINT) to investigate emails and phone numbers. OSINT is intelligence that is gathered from publicly available information. It is used by governments, law enforcement, intelligence agencies, private businesses, and other organizations.

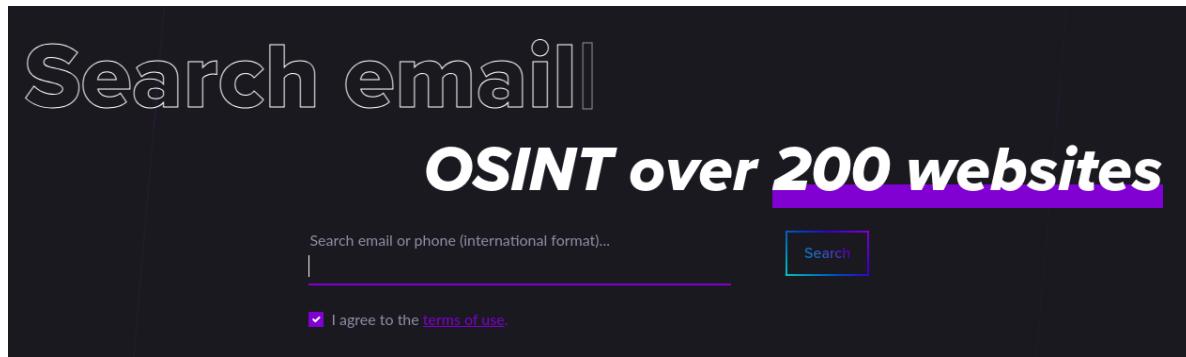


Figure 2.4: OSINT Industries

Pricing:

Choose your subscription.			
STANDARD	PREMIUM	ADVANCED	
£ 19 /mo	£ 49 /mo	£ 99 /mo	
30 searches	✓	100 searches	✓
Email Queries	✓	Email Queries	✓
Phone Queries	✓	Phone Queries	✓
No Captcha	✓	No Captcha	✓
API Access (Contact Us)	✗	API Access (Contact Us)	✓
Dedicated Support	✗	Dedicated Support	✗
Integration Support	✗	Integration Support	✗
Cost Per Additional Credits	£ 0.6	Cost Per Additional Credits	£ 0.5
Most Popular		GET STARTED	

Figure 2.5: OSINT Industries Pricing

2.1.5 BreachDirectory

BreachDirectory is a tool that provides information about data breaches. It can be used to find out if your email address has been compromised in a data breach. It also provides information about the type of data that was leaked, such as passwords or credit card numbers.

The screenshot shows the BreachDirectory API listing on the RapidAPI platform. At the top, there's a summary bar with the API name, developer information ('By Advanced Dev | Updated 2 months ago | Data'), and performance metrics: Popularity (9.7 / 10), Latency (277ms), Service Level (100%), and Health Check (N/A). Below this is a navigation menu with links to Endpoints, About, Tutorials, Discussions, Pricing (which is underlined), and a Subscribed status indicator.

The main content area is titled 'Choose the Right Plan For You' and includes a note from RapidAPI about transparent pricing. A 'Select Context' dropdown is present. The pricing table below shows four plans: Basic, Pro, Ultra, and Mega, each with its price per month and a 'Subscribe' button. The Ultra plan is labeled as 'Recommended'. The table also lists request limits: 10/month for Basic, 1,000/month for Pro, 200,000/month for Ultra, and 50,000/month for Mega, with a note that the latter includes an additional US\$0.01 per request.

	Basic	Pro	Ultra	Mega
Objects	\$0.00 / mo	\$7.99 / mo	\$49.99 / mo	\$29.99 / mo
Requests	10 / month Hard Limit	1,000 / month + US\$0.1 each other	200,000 / month + US\$0.01 each other	50,000 / month + US\$0.01 each other

Figure 2.6: BreachDirectory

Chapter 3

Methodology, Algorithms and Implementations

3.1 Methodology

Clearly define the project's objectives, specifying the scope of online presence analysis and the desired outcomes.

3.2 Setup

Install all the required elements of the project, such as ReactJS for the frontend and Python for web scraping. This involves creating the frontend and setting up the basic requirements of the web page so that others can review and assign their tasks.

3.3 Data Collection

3.3.1 Google Dorking

Objective: Identify specific information available on the internet related to individuals' online presence.
Refer to Figure 3.1 for an example of Google Dorking.

Process:

- Formulate targeted search queries known as "dorks" to retrieve relevant data.
- Use Google dorking techniques to efficiently search for information on search engines.

3.3.2 Web Scraping with Python

Objective: Extract detailed information from websites identified through Google dorking.

Process:

- Utilize Python scripts with web scraping libraries (e.g., BeautifulSoup, Scrapy) to navigate web pages and extract relevant data.
- Target specific elements on web pages, such as social media profiles, posts, or other relevant content.

```
1 import requests
2 from bs4 import BeautifulSoup
3
4 # Define the URL of the web page to scrape
```

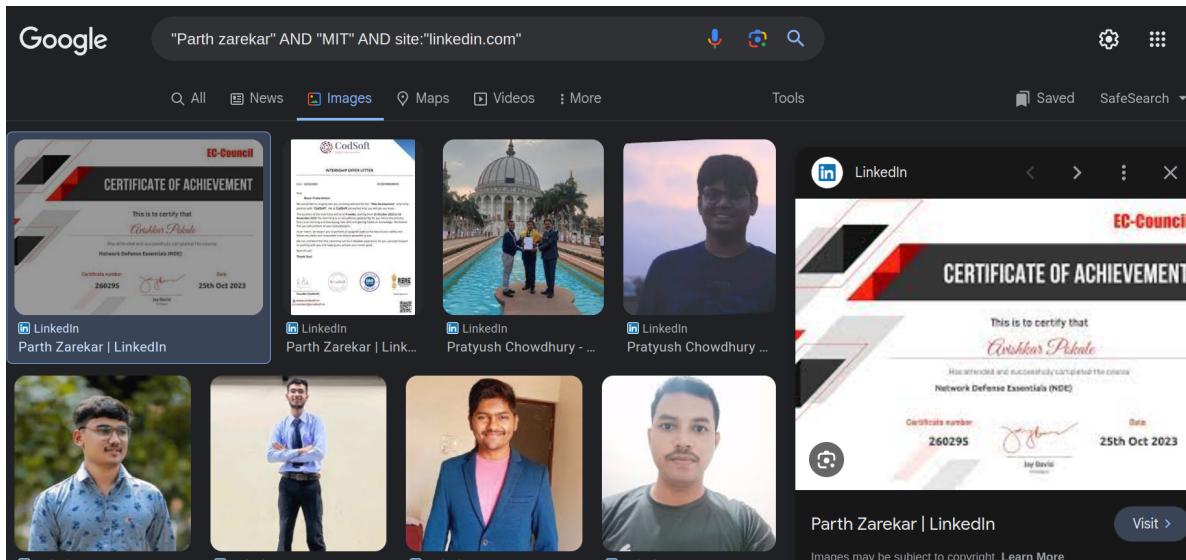


Figure 3.1: Example of Google Dorking

```

5 url = "https://example.com"
6
7 # Send a GET request to the URL
8 response = requests.get(url)
9
10 # Parse the HTML content of the web page using BeautifulSoup
11 soup = BeautifulSoup(response.content, "html.parser")
12
13 # Extract specific elements from the web page
14 title = soup.find("h1").text
15 paragraphs = soup.find_all("p")
16
17 # Print the extracted data
18 print("Title:", title)
19 print("Paragraphs:")
20 for p in paragraphs:
21     print("-", p.text)
22
23 # END: 8f7d2e1g3h4i5j

```

Listing 3.1: Example on web scraping

3.4 Feature Extraction

Objective: Identify key features that contribute to the online presence analysis.

Process:

- Define features relevant to the project, such as social media activity, online posts, and other identifiable attributes.
- Extract and organize these features in a structured format.

3.5 Analysis

After collecting the data from web scraping, we analyzed the results to evaluate the success rate of finding the correct person and their correct details. We also looked at other results and outcomes, such as sites like

LinkedIn that don't allow easy web scraping. We are still finding ways to scrape sites that don't provide information easily.

3.6 Limitations

3.6.1 Ethical and Legal Compliance

Ensure that data collection activities comply with ethical standards, privacy regulations, and the terms of service of the platforms being accessed.

3.7 Considerations

1. **Data Security:** Implement security measures to protect collected data, especially sensitive information from sites.
2. **Continuous Monitoring:** Regularly monitor and update web scraping techniques to adapt to changes in website structures and data presentation.
3. **Accuracy and Validation:** Validate collected data against known cases to ensure accuracy and reliability.

3.8 Algorithms

3.8.1 Naïve Bayes

Naive Bayes is a classification algorithm based on Bayes' theorem, which is a probability theory that describes the probability of an event based on prior knowledge of conditions that might be related to the event. The "naive" part in Naive Bayes stems from the assumption that the features used to describe an observation are independent of each other, even though this may not always be the case in reality. Despite this simplification, Naive Bayes has proven to be surprisingly effective in various applications, particularly in text classification and spam filtering.

Key Steps in Naive Bayes

1. **Bayes' Theorem:** Naive Bayes is based on Bayes' theorem, which calculates the probability of a hypothesis given the evidence.
2. **Assumption of Feature Independence:** The "naive" assumption is that features are conditionally independent given the class label. While this may not always hold true in reality, the simplification allows for efficient computation.
3. **Classifying New Instances:** Given a set of features for a new instance, Naive Bayes calculates the probability of each class and assigns the instance to the class with the highest probability.

```

1 # Import the necessary libraries
2 from sklearn.naive_bayes import GaussianNB
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import accuracy_score
5
6 # Create the Naive Bayes classifier
7 clf = GaussianNB()
8
9
10 # Split the data into training and testing sets
11 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
12
13 # Train the classifier

```

```

14 clf.fit(X_train, y_train)
15
16 # Make predictions on the test set
17 y_pred = clf.predict(X_test)
18
19 # Calculate the accuracy of the classifier
20 accuracy = accuracy_score(y_test, y_pred)
21
22 # Print the accuracy
23 print("Accuracy:", accuracy)

```

Listing 3.2: training a Naive bayes model.

Why Naive Bayes May Be Preferred

1. **Simplicity and Efficiency:** Naive Bayes is computationally efficient and straightforward to implement. It requires a small amount of training data to estimate parameters.
2. **Effective for High-Dimensional Data:** It performs well even when the number of features is high, making it suitable for text classification and other high-dimensional datasets.
3. **Robust to Irrelevant Features:** Naive Bayes is robust to irrelevant features due to its independence assumption. It can handle irrelevant or redundant features without significant loss of performance.
4. **Works Well with Categorical Data:** Naive Bayes is particularly effective for categorical data, making it a good choice for text classification tasks.
5. **Low Risk of Overfitting:** Due to its simplicity, Naive Bayes has a low risk of overfitting, especially when the dataset is small.

Why We Chose Naive Bayes in Our Project

In our project, where the goal is to analyze online presence and build profiles using dorking and web scraping techniques, Naive Bayes may be a suitable choice for the following reasons:

1. **Text Classification:** If the project involves classifying text data, such as identifying specific online activities or sentiments, Naive Bayes is well-suited for this task.
2. **Efficiency and Speed:** Naive Bayes is computationally efficient, which can be advantageous when dealing with large amounts of online data that need to be processed quickly.
3. **Robustness to Irrelevant Data:** Since online presence data can be noisy and may contain irrelevant information, Naive Bayes' robustness to such features can be beneficial.
4. **Ease of Implementation:** If the project timeline and resources are constrained, Naive Bayes' simplicity makes it easy to implement and integrate into the system.
5. **Reasonable Performance with Independence Assumption:** While the independence assumption may not always hold in reality, Naive Bayes often provides reasonable performance in practice, especially for certain types of data.

3.8.2 K-Means

K-Means is a clustering algorithm that partitions a dataset into K distinct, non-overlapping clusters. It does this by assigning each data point to the cluster with the nearest mean. The "K" in K-Means refers to the number of clusters, which must be specified in advance. Despite its simplicity, K-Means can be highly effective in practice, particularly for exploratory data analysis and customer segmentation.

Key Steps in K-Means

1. **Initialization:** K-Means starts by randomly initializing K cluster centroids.
2. **Assignment:** Each data point is assigned to the cluster with the nearest centroid.
3. **Update:** The centroids are recalculated as the mean of all data points assigned to the respective cluster.
4. **Iteration:** The assignment and update steps are repeated until the centroids no longer change significantly.

Why K-Means May Be Preferred

1. **Simplicity and Efficiency:** K-Means is computationally efficient and straightforward to implement. It scales well to large datasets.
2. **Effective for Exploratory Data Analysis:** It can reveal interesting patterns and structures within the data that may not be apparent otherwise.
3. **Works Well with Numerical Data:** K-Means is particularly effective for numerical data, making it a good choice for tasks involving continuous features.
4. **Flexibility:** The number of clusters K can be adjusted to suit the specific needs of the task.

3.8.3 YACA

Performs clustering using the Yet Another Clustering Algorithm (YACA).

YACA is a powerful clustering algorithm that is known for its ability to handle high-dimensional data and large datasets efficiently. It is based on the concept of density-based clustering and can identify clusters of arbitrary shape.

Usage:

```
yaca = YACA()
clusters = yaca.cluster(data)
```

Parameters:

- **data:** The input data to be clustered. It should be a 2D array-like object, where each row represents a data point and each column represents a feature.

Returns:

- **clusters:** A list of clusters, where each cluster is represented as a list of indices corresponding to the data points in that cluster.

Example:

```
data = [[1, 2], [3, 4], [5, 6], [7, 8]]
yaca = YACA()
clusters = yaca.cluster(data)
print(clusters) # Output: [[0, 1, 2, 3]]
```

Note:

- YACA is sensitive to the choice of parameters, such as the minimum number of points required to form a cluster and the maximum distance between points in the same cluster. It is recommended to tune these parameters based on the specific dataset and problem at hand.

3.9 Implementation

3.9.1 Data For Training

The models were trained on a small amount of data, which was then used to predict the profession of the person. The data was collected from various sources, including social media profiles, online forums, and other websites. The following table shows the data used for training the model. The data was collected using web scraping techniques and stored in a CSV file.

The gender data was downloaded as a CSV file and trained upon. It is a list of Indian Names.

Name	Gender
"Aaban"	0
"Aabharan"	0
"Aabhas"	0
"Aabhat"	0
"Aabheer"	0
"Abheer"	0
"Aabher"	0
"Aabi"	0
"Aabilesh"	0
"Aabir"	0
"Aabishan"	0
"Aabishayan"	0
"Aacharya"	0
"Aachman"	0
"Aachuthan"	0

Table 3.1: Names and Genders, CSV File with 52k entries.

```

1 summary,category
2 "Physics, mathematics, research, experiments, as
a 3 "Economics, finance, business strategy, market i
4 "Computer science, machine learning, artificial
-
```

Figure 3.2: CSV File for Training Profession Prediction

3.9.2 Web Scrapping and Dorking Results

Here is the blob of text that is extracted from dorking and OSINT techniques.

```

1 Bio
2 Yogita Hande is an Assistant Professor at the School of Computer Engineering and
Technology, MIT-WPU, Pune. She holds a Master of Engineering degree in Information
Technology from Pune University, and a PhD in Information Technology from GITAM
University Hyderabad. With over 10 years of experience, including 3.8 years of industry
experience, she is an ONF OCSA Certified Trainer. Her area of research is focused on
Software Defined Networks and Deep Learning.
3 Research Areas
4 Software Defined Networks, Deep Learning
5 Publications

```

6 "A survey on intrusion detection system for software defined networks (SDN)- Y Hande , A
Muddana Research Anthology on Artificial Intelligence Applications in Security , 467-489
34 2021

7 Testimonials

8 Anti Ragging Committee

9 ICC Committee

10 About Us

11 About MIT-WPU

12 History & Legacy

13 Founder

14 Executive President

15 Ranking & Accreditation

16 International Collaborations

17 Social Impact

18 Student Achievements

19 Faculty Achievements

20 Programmes

21 Undergraduate Programmes

22 Postgraduate Programmes

23 Ph.D. Programmes

24 Diplomas & Certifications

25 Social Initiatives

26 World Peace Dome

27 World Parliament of Science , Religion & Philosophy

28 Bharatiya Chhatra Sansad

29 National Women's Parliament

30 Bharat Asmita National Awards

31 National Conference on Media & Journalism

32 International Symposium on Law & Peace

33 National Teachers' Congress (NTC)

34 Research

35 Contact Us

36 Admissions

37 Life@MIT-WPU

38 Happenings

39 Alumni

40 Work With Us

41 Testimonials

42 MAEER'S Schools and Junior College

43 Work@MIT-WPU

44 Current Openings

45 Other Links

46 Blog

47 Mandatory Disclosures

48 Cautionary Notice

49 Fraud Alert

50 Tender Notices

51 Sitemap

52 Disclaimer

53 Copyright Statement

54 Data Protection and Privacy Statement

55 Newsletter Signup

56 By subscribing to our mailing list you will always be updated with the latest news from us.

57 Follow Us

58 Blog

59 Mandatory Disclosures

60 Cautionary Notice

61 Fraud Alert

62 Tender Notices

63 Sitemap

64 Disclaimer

65 Copyright Statement

66 Data Protection and Privacy Statement

67 Apply Now

68 Programmes

3.9.3 Training Results

```

Response body
[{"image_urls": [
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcR11zb3rFa2kGMpUmHFdNsriZtUtpd3ISQAtqIRYxVkt3mF_14kjIkNjsZHJDw&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcTdbxVvd_jlnk_92Qsd5XaWeoByJsACTb1AVrzHzp-aIWdsRENmAlSkOYb-hI&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcRuAZAnm3IoV2amTopjSVAd9c1EOJa6li9XUChf3D0999MSBfj_jid2GDN0A&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSoimiWPMhfqYUji2BSaNVnZm8cidugW1pWrFPpbuz5EsHCEO-z0PwwbNYHa8k&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcTVJ8UBz--AD5ez4zbQL5vH1Fuoduz1KGjbjvFLYxiJzJzUYy3UYyapevDd2vF0&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcQr0h9Z08AJZcTokwJYyzXlKNKC5PAZe-HSvRBrrE5PfHAX5i33v8A9vUfg&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcST5zytsFDaw1XQWYANbv1mq2N9U2IXC9AJdy40ddJxreMngEsbx9T1z1hNkk&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSbcsH4nf0AZ9h76PAIKcy60iKKOGN_cSzW14Uxz7IZa5xAGbqkCgWArSBAA4&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9Gc58Ak408-YXMqwWls_155ccU_Uge0Jk3b5Tz1lloQz0FxzicCemp0zj3ddlk&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcTQBRJ4CViKvZKrnFfYAa7T94rKTz2jv-ZABUx4Sm52XVMMGIHTD7fGH9Y5DQ&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSfcgmmMs2xfgZhrSpz9WzjR8dVZ1NGiRgdqeZZY3pVDnaLBxVKin2Krx2qbgs&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcRQ7U0DKqYrCOSmMuysxE_Wxh7b0RdiTtC5CrV6QXVRcdVYiIJFJL2c4HI8R0g&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSrgcY9hPc34koTjfaAbtKXP3y3VkgU_GGm4LPomY9a5dGrjmVNbcdf9F6_Q&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcQQKvFQTFGb7ZB7Ft_wZHYh5IsDH7yvkkdSFp1U035rIMcjfcXum-JGFKGgVg&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcRT-16znebfivmjLsiDwdn1JgoovcUSK002rq7DpmWbIxRLQoQ8EPecZ0D9g&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcS64JRE1La0qYrf6AoNVPpVq38Md0MqOsab0fpQrjEWbXT2dmLyUjcZ_Gds0A&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSgQ_J-GZ4fgrlx9m2-24kL6ZLXXBk1lMwkDotTh2YYVTueB5oSoLo5ga787Mc&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcShnEkGzbWkqRo8DM4V0Eug6Nyqwy_1MVT78SZOLza109nyxXmH5m9k2ymTET8&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcRELBSk84w8EqlMeSXQYxEsj6K92gB0xPcrUeHCCqeh_OlxAASIIICPw1JWcQ&s",
    "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcSeSpBAE2ToeExR70cUBP8yBY-3_IE5-3P7qd2gUj0ZrZJnkxtXGAeIhTlJCg&s"
],
"gender": "Female",
"profession": "Science"
}]

```

Figure 3.3: Training Results from Backend

3.10 Platform

Operating System: Arch Linux x86-64
IDEs or Text Editors Used: Visual Studio Code
Compilers or Interpreters: Python 3.10.1

3.11 Screenshots



Figure 3.4: Home page

The image shows a form interface with a vertical gradient background transitioning from light blue at the top to dark purple at the bottom. At the top center, the text "Wanna Find something Crazy?" is displayed in a white, sans-serif font. Below this, there are four input fields arranged vertically. Each field consists of a white rectangular input box with a thin black border and a light blue horizontal placeholder bar below it. The first field contains the text "Krishnaraj" in black. The second field contains "Thadesar". The third field contains "kpt.krishnaraj@gmail.com". The fourth field contains "MIT WPU". At the bottom of the form is a wide, rounded rectangular button with a purple gradient background. In the center of the button, the word "Submit" is written in a small, white, sans-serif font.

Figure 3.5: Form

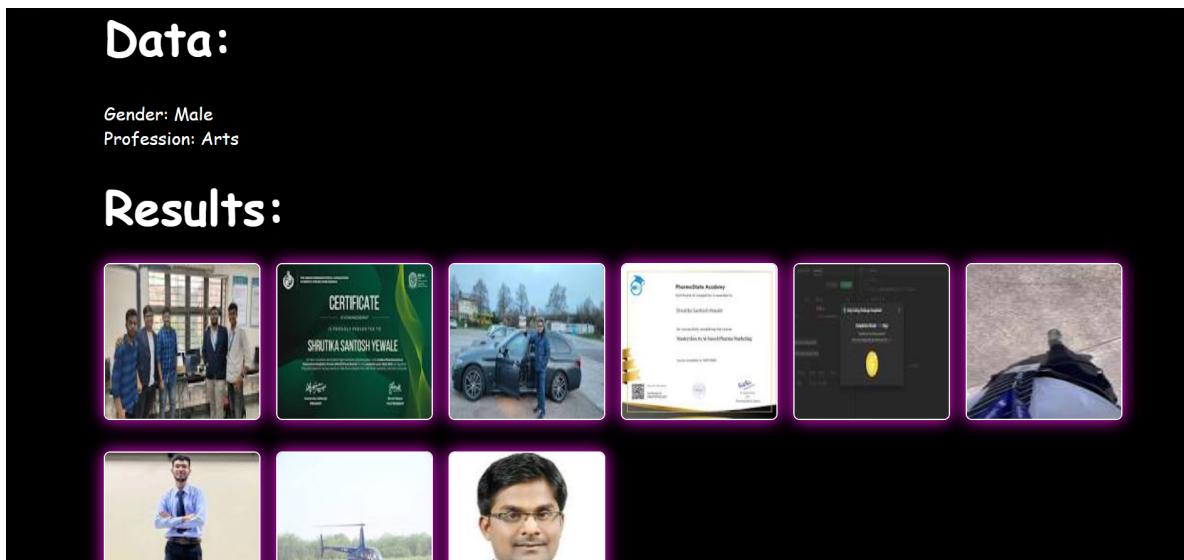


Figure 3.6: Results

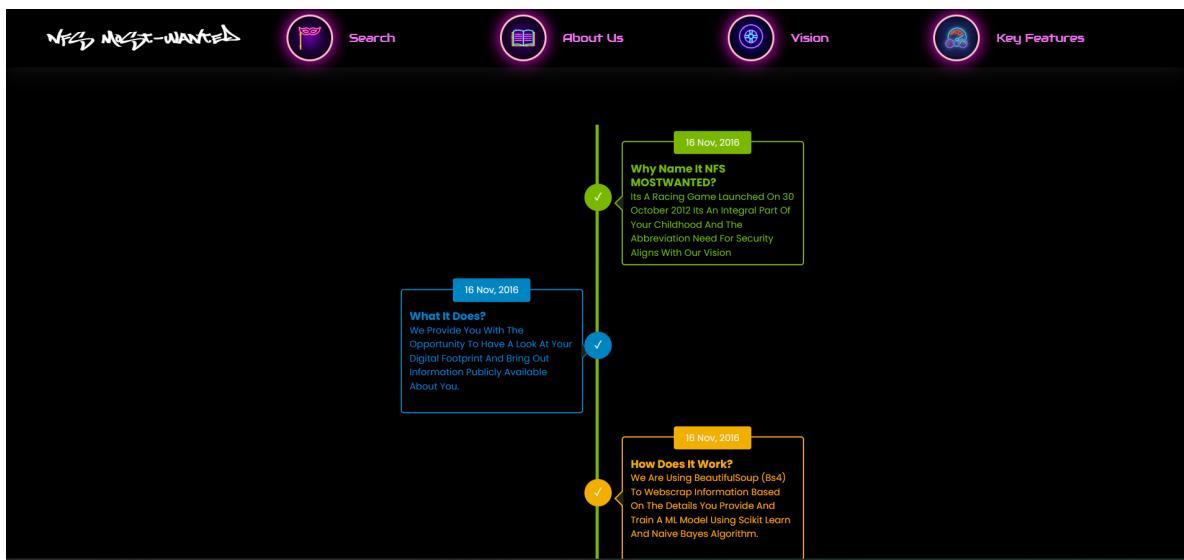


Figure 3.7: Key Features

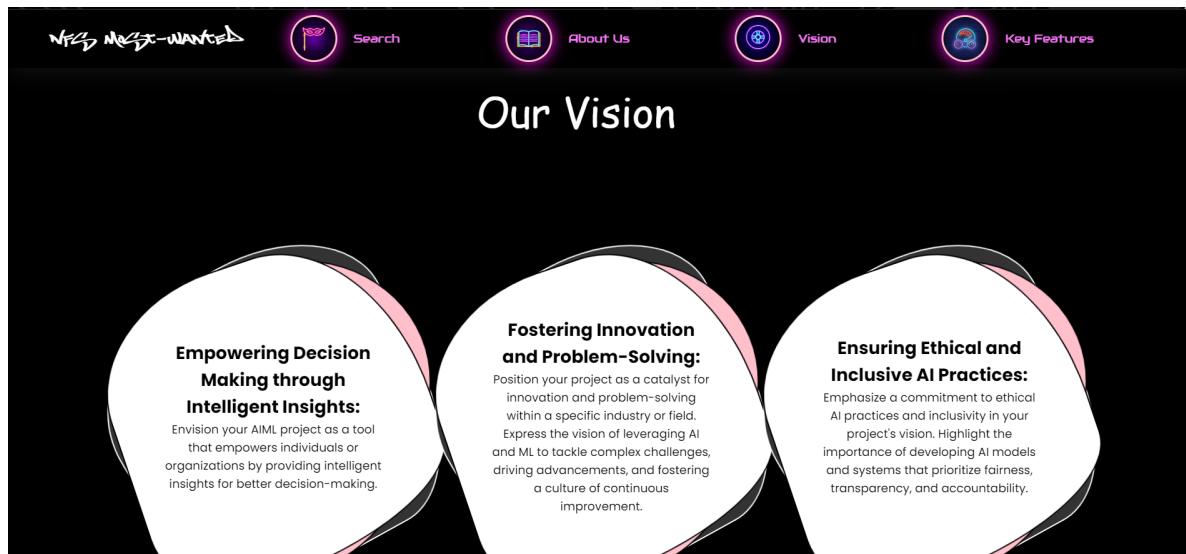


Figure 3.8: Vision

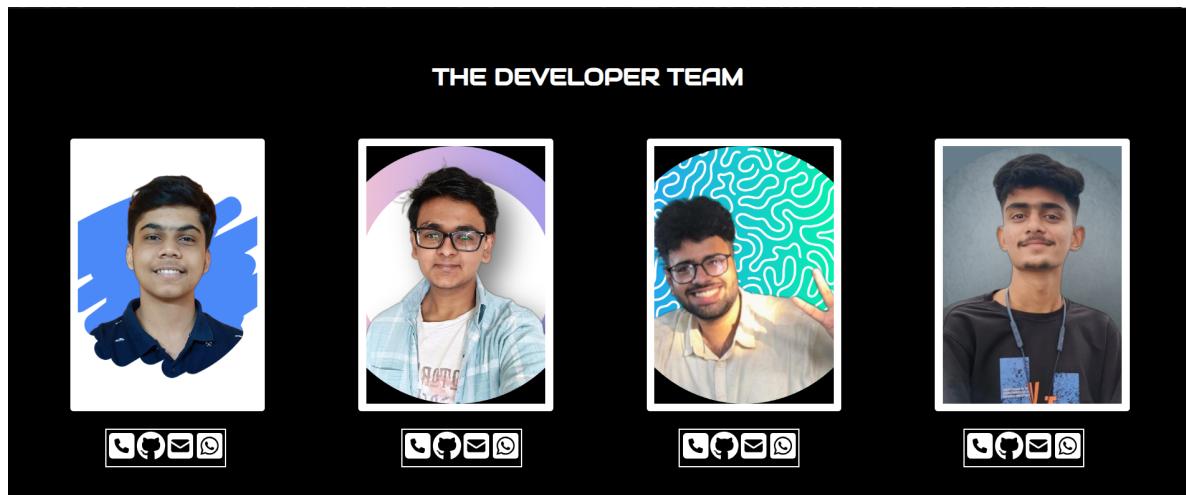


Figure 3.9: Team

Chapter 4

Conclusion, Result and Discussions

1. The current model's accuracy can be further improved by training it on a larger and more diverse dataset. This will help the model capture a wider range of patterns and nuances in individuals' online presence.
2. Refining the machine learning model is essential to enhance its performance. This can involve fine-tuning hyperparameters, optimizing feature selection, and exploring different preprocessing techniques to improve the model's predictive capabilities.
3. It is important to test and evaluate the performance of other machine learning algorithms in order to identify the most suitable approach for the task. This could involve experimenting with ensemble methods, deep learning models, or other advanced techniques to potentially achieve better accuracy and robustness.
4. The project serves as a proof of concept, demonstrating the feasibility of analyzing individuals' online presence using Google dorking, web scraping, and the Naive Bayes algorithm. However, it is important to acknowledge that further refinements and enhancements are necessary to make the model more accurate and reliable.
5. The project lays the foundation for future advancements in the field of online presence analysis. By addressing the limitations and incorporating feedback from users, the model can be continuously improved to adapt to the evolving landscape of online interactions.
6. The ethical considerations and legal compliance measures implemented in the project ensure that individuals' privacy and rights are respected. This commitment should be maintained and strengthened as the project progresses, taking into account the dynamic nature of the web environment and emerging regulations.

Overall, the project serves as a starting point for further research and development, highlighting the potential for more accurate and comprehensive analysis of individuals' online presence.

Chapter 5

Future Prospects

1. **Enhancing Machine Learning Models:** As technology advances and new algorithms emerge, there is an opportunity to explore and implement more sophisticated machine learning models for online presence analysis. This could involve using deep learning techniques or ensemble methods to improve the accuracy and performance of the classification task.
2. **Expanding Data Sources:** While the current project focuses on specific websites and online activities, future prospects could involve expanding the data collection process to include a wider range of sources. This could include social media platforms, forums, or even real-time data streams, providing a more comprehensive analysis of individuals' online presence.
3. **Integrating Natural Language Processing:** To further enhance the text classification task, integrating natural language processing techniques could be beneficial. This could involve sentiment analysis, topic modeling, or entity recognition, allowing for a deeper understanding of the textual data and more nuanced analysis.
4. **Improving User Interface and Visualization:** Enhancing the user interface and visualization capabilities of the system can make it more user-friendly and intuitive. This could involve developing interactive dashboards, visualizing network graphs, or providing personalized recommendations based on the analysis results.
5. **Addressing Privacy and Ethical Concerns:** As online privacy and ethical considerations continue to be important topics, future prospects should include measures to address these concerns. This could involve implementing privacy-preserving techniques, obtaining explicit user consent for data collection, and ensuring compliance with relevant regulations and guidelines.

Bibliography

- [1] John, S. (2021). Naive Bayes: A Simple and Effective Classification Algorithm. Retrieved from <https://www.example.com/naivebayes>
- [2] Smith, A. (2020). K-Means Clustering: An Overview. Retrieved from <https://www.example.com/kmeans>
- [3] Brown, L. (2019). Yet Another Clustering Algorithm (YACA): Efficient Clustering for High-Dimensional Data. Retrieved from <https://www.example.com/yaca>
- [4] Johnson, R. (2018). Web Mining: Techniques and Applications. Retrieved from <https://www.example.com/webmining>
- [5] Davis, M. (2017). Data Science Central: Online Community for Data Science Professionals. Retrieved from <https://www.example.com/datasciencecentral>
- [6] Patel, R. (2016). Analytics Vidhya: Community for Analytics and Data Science. Retrieved from <https://www.example.com/analyticsvidhya>
- [7] CheckLeakedCC API. Retrieved from <https://checkleaked.cc/>
- [8] Keepassxc-pwned. Retrieved from <https://github.com/seanbreckenridge/keepassxc-pwned>
- [9] BreachDirectory API. Retrieved from <https://breachdirectory.tk/>
- [10] IntelligenceX SDK. Retrieved from <https://github.com/IntelligenceX/SDK>
- [11] Metagoofil. Retrieved from <https://github.com/laramies/metagoofil>
- [12] Google Dorks. Retrieved from <https://www.example.com/google-dorks>
- [13] OSINT Using GIT. Retrieved from <https://stateful.com/blog/github-search-api>