

## Bytexl's guided project

### Final Project report

Name of the educator	Mounesh Gouda
Project title	Ai Based Fake Job Post Prediction
Tools / platforms used	Local System,NLP, MachineLarning,Colab

**Title:** Ai Based Fake Job Post Prediction

#### About the Project :

The rise of online job portals and social media has opened new opportunities for job seekers but has also led to an increase in fraudulent job postings. These fake job posts often mislead candidates, leading to scams, identity theft, or financial loss. The goal of this project is to develop an AI-based system that can effectively predict and identify fake job posts from legitimate ones. By leveraging machine learning and natural language processing (NLP) techniques, this system can analyze the content of job listings and detect patterns that indicate fraudulent activity.

#### System Requirements:

- **Software:**
  - Programming Language: Python 3.8+
  - Key Libraries:
    - Machine Learning: Scikit-learn, TensorFlow/Keras, XGBoost
    - NLP: NLTK, SpaCy, Hugging Face Transformers
    - Web Scraping: BeautifulSoup, Selenium
    - Web Framework: Flask (for API and deployment)
    - Development Tools: Visual Studio Code or PyCharm, Git for version control
- **Hardware:**
  - Processor (CPU): Intel Core i5 or higher (Recommended: i7 or better for faster model training)
  - RAM: Minimum 8 GB (Recommended 16 GB+)
  - Storage: Minimum 100 GB (SSD recommended for faster data handling)
  - GPU (Optional): NVIDIA GTX 1060 or better (Recommended for deep learning tasks)

#### Functional Requirements:

- Job Post Submission & Analysis: Users can submit job posts for classification (Fake or Legitimate) using text input or bulk upload.
- Fraud Detection: The system identifies suspicious patterns (e.g., unrealistic salaries, vague descriptions) and validates company info.
- Real-Time Alerts & Reports: Users receive notifications and detailed reports on flagged job posts, with reasons for classification.
- Continuous Improvement & Data Privacy: The system continuously retrains with new data to improve accuracy, while ensuring compliance with data privacy regulations.

## User Interface Requirements:

Below is the user interface flow, with input fields and output sections, based on the steps you've described:

### Inputs and Outputs:

- **Inputs:** Username and Password for login. CSV file containing job post data for training. Job post description to predict if it's fake or legitimate. Click on "Predict" button to analyze the job post. Click on "Try Another Post" to reset the input field for a new job post.
- **Outputs:** Successful login or error message for incorrect credentials. Success message and model training initiation or error message if upload fails. Prediction result (Fake/Legitimate) with confidence score and reasoning. Classification result (Fake/Legitimate) with confidence score and highlighted suspicious keywords. Cleared input field ready for a new job post description.

### List of Subsystems:

1. **User Authentication Subsystem:** login, registration, and session manage
2. **Data Upload and Preprocessing Subsystem:** Manages CSV file upload, data validation, and preprocessing for model training.
3. **Model Training Subsystem:** Trains the AI model on the uploaded data and evaluates its performance.
4. **Prediction Subsystem:** Classifies job posts as Fake or Legitimate based on user inputs using the trained model.
5. **User Interface Subsystem:** Manages the frontend interface for job post submission, file upload, and displaying prediction results.

### Applications in Other Contexts:

AI-based systems used for Fake News Detection, Spam Email Filtering, Fraudulent Account Detection, Phishing Website Detection, and Product Review Authenticity all rely on similar machine learning and natural language processing techniques to analyze patterns, identify anomalies, and improve security and trust across various digital platforms.

### Designing of Test Cases:

#### ➤ Test Cases for Data Upload and Preprocessing

##### 1. Valid CSV Upload:

**Input:** Correctly formatted CSV.

**Expected Output:** Data uploaded and parsed without errors.

##### 2. Invalid CSV Format:

**Input:** Malformed CSV (e.g., missing headers).

**Expected Output:** Error message indicating invalid format.

##### 3. Missing Data:

**Input:** CSV with missing job description or title.

**Expected Output:** Missing fields flagged for correction.

➤ **Test Cases for Model Prediction**

**1.Fake Job Prediction:**

**Input:** Known fake job post.

**Expected Output:** Model predicts "Fake".

**2.Real Job Prediction:**

**Input:** Real job post.

**Expected Output:** Model predicts "Real".

**3.Ambiguous Job Post:**

**Input:** Mixed job post (real and fake elements).

**Expected Output:** Model predicts "Uncertain" or gives a confidence score.

➤ **Test Cases for UI/UX**

**1.User Login:**

**Input:** Valid username and password.

**Expected Output:** Successful login and redirection to the home page.

**2.File Upload Button:**

**Input:** User clicks upload and selects CSV.

**Expected Output:** File is successfully uploaded and displayed.

**3.Invalid File Upload:**

**Input:** Non-CSV file selected.

**Expected Output:** Error message: "Invalid file type."

➤ **Test Cases for Performance and Load Testing**

**1.Large Dataset Processing:**

**Input:** CSV with 10,000+ job posts.

**Expected Output:** File processed within 5 seconds.

**2.Real-Time Prediction Speed:**

**Input:** Job post for prediction.

**Expected Output:** Prediction made in under 2 seconds.

**3.Concurrent Users:**

**Input:** Multiple users (5-10) using the system.

**Expected Output:** No performance issues under load.

### **Future Work:**

The future work includes continuously retraining the model with new data to enhance prediction accuracy, integrating the system with popular job platforms for real-time fake job post detection, extending multilingual support to detect fake job posts in various languages, adding advanced fraud detection features like behavioral analysis to identify suspicious patterns, and implementing a user feedback mechanism to improve the system by allowing users to report false positives and negatives.

### **References:**

1. Yang, D. D. K. G. F., & Yu, L. M. (2020). Online Recruitment Fraud: An Analysis of Job Posting Scams. *Journal of Cybersecurity and Privacy*, 4(2), 123-145.
2. Smith, A., & Johnson, J. (2021). Detecting Fake Job Postings with Machine Learning. *International Journal of Information Technology*, 12(3), 101-114.
3. Chen, M., & Kumar, R. (2022). Natural Language Processing for Fraud Detection. *Computational Linguistics*, 48(1), 67-85.
4. Lee, K., & Park, S. (2021). Feature Engineering for Job Posting Classification. *Data Mining and Knowledge Discovery*, 35(4), 789-806.
5. Wang, J., & Zhang, T. (2022). Evaluating Machine Learning Models for Online Job Fraud Detection. *Journal of Machine Learning Research*, 23(1), 200-220.
6. Thompson, F. G., & Singh, N. L. (2021). Combining Textual and Structured Data for Job Posting Verification. *Artificial Intelligence Review*, 54(2), 345-367.
7. Patel, R. B., & Verma, A. S. (2022). Deep Learning Approaches for Fake Job Detection. *Journal of Neural Computing and Applications*, 34(5), 1010-1025.
8. Kim, H. J., & Nguyen, L. T. (2023). Real-Time Detection of Job Posting Scams. *IEEE Transactions on Information Forensics and Security*, 18(3), 512-529.
9. Smith, M. R., & Johnson, T. (2020). Understanding the Impact of Fake Job Postings on Job Seekers. *Journal of Business Ethics*, 165(3), 456-472.
10. Brown, A. F., & Wilson, K. P. (2021). Legal Implications of Job Posting Fraud. *Journal of Internet Law*, 25(4), 45-60.
11. Garcia, L. T., & Thompson, M. E. (2022). Predictive Analytics for Online Job Fraud Detection. *International Journal of Data Science and Analytics*, 6(2), 89-104.
12. Jones, S. P., & Roberts, A. C. (2021). Social Network Analysis for Job Posting Verification. *Social Network Analysis and Mining*, 11(1), 78-92.

## **Reflection of the Project Creation:**

### **Technical Challenges Encountered**

- Data Preprocessing:Organizing unstructured job post data for model training was time-consuming.
- Model Training:Overcoming overfitting and selecting relevant features was challenging due to an imbalanced dataset.
- Real-Time Prediction:Ensuring fast predictions without sacrificing accuracy, especially for long job descriptions, required optimization.

### **How Existing Software Engineering Knowledge Helped**

Data Handling:My knowledge of data structures and algorithms helped with efficient data processing and feature extraction.

Version Control: Git enabled smooth collaboration and code tracking across the project.

UI/UX:Applying UI/UX principles made the interface intuitive and user-friendly for non-technical users

### **Benefits Experienced While Working on This Project**

Web Scraping Skills: I gained practical experience in web scraping, which helped in collecting real-time job post data from different sources for training the model.