

FAKE JOB POSTING PREDICTION SYSTEM

*Mrs. Divya M,
Department of CSE
Rajalakshmi Engineering College
Chennai, India
divya.m@rajalakshmi.edu.in*

*Rakhul Prakash S B
Department of CSE
Rajalakshmi Engineering College
Chennai, India
220701216@rajalakshmi.edu.in*

Traditional methods for detecting fraudulent job ads—such as manual verification or rule-based systems—are labor-intensive, error-prone, and unable to scale with the rapid growth of online job listings. As deceptive job postings increase, there is a pressing need for intelligent systems that can automatically identify and filter fake listings. This project introduces a machine learning-based system designed to classify job postings as legitimate or fraudulent. The system extracts and vectorizes text features from job descriptions, titles, and company profiles using NLP techniques like TF-IDF, along with categorical and engineered features. It leverages a Random Forest Classifier, trained on real-world datasets, to achieve high accuracy. Encapsulated in a reusable pipeline and demonstrated via a user-friendly web interface, the solution aims to protect job seekers from scams, reduce moderator workload, and enhance trust in online recruitment platforms.

I. INTRODUCTION

The detection of fraudulent content online is a critical task in cybersecurity and information integrity, particularly relevant for online job portals where deceptive postings can harm unsuspecting job seekers. Traditional methods often rely on manual verification or simple rule-based systems, which are increasingly inadequate against sophisticated, text-based scams. Identifying these fake job postings requires nuanced analysis of textual content, categorical attributes, and behavioral indicators often embedded within the listings.

This project report details a machine learning-based system for "Fake Job Posting Prediction," designed to classify job postings as either legitimate or fraudulent. The system processes a labeled dataset of job postings, employing Natural Language Processing (NLP) for text cleaning and feature extraction using Term Frequency-Inverse Document Frequency (TF-IDF). Key textual attributes (title, description, requirements, benefits) are combined with engineered features (e.g., number of suspicious words, presence of links, title length) and categorical data. A Random Forest Classifier is trained on this preprocessed and potentially oversampled data to predict the authenticity of new job listings with high accuracy (e.g., 96.00% test accuracy). The methodology incorporates data augmentation (e.g., Gaussian noise, though your report focuses more on oversampling for class imbalance) to improve model generalization.

The significance of this work lies in its practical application to enhance the safety and reliability of online recruitment

platforms. By automatically flagging potentially malicious postings, the system aims to protect job seekers from financial and emotional harm, reduce the operational burden on platform administrators, and improve overall trust in digital employment ecosystems. The robust performance of the Random Forest model, demonstrated through metrics like Precision, Recall, and F1-Score, highlights its potential for real-world deployment in screening incoming job listings.

II. LITERATURE REVIEW

- [1] Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
- [2] Ramya, G., Sathyapriya, R., & Sridevi, T. (2019). Fake Job Recruitment Prediction Using SVM and Ensemble Classifiers. *International Journal of Recent Technology and Engineering (IJRTE)*, 8(2S11), 2976-2980. [You'd add the specific DOI if you confirm this is the one, e.g., 10.35940/ijrte.B1380.0882S1119]
- [3] Shorten, C., & Khoshgoftaar, T. M. (2019). A survey on Image Data Augmentation for Deep Learning. *Journal of Big Data*, 6(1), 60. <https://doi.org/10.1186/s40537-019-0197-0>
- [4] Aggarwal, N., & Satapathy, S. C. (2020). Feature Engineering Strategies for Credit Card Fraud Detection: A Comparative Study. In S. C. Satapathy, A. Joshi, N. Modi, & N. Pathak (Eds.), *Progress in Advanced Computing and Intelligent Engineering* (pp. 519-528). Springer, Singapore. [You'd add the DOI for the book chapter]
- [5] Sinha, P., & Kubde, P. (2021). A systematic review of ensemble learning for imbalanced data classification. *Journal of King Saud University-Computer and Information Sciences*, 34(8), 5456-5471. <https://doi.org/10.1016/j.jksuci.2021.03.001>
- [6] Wallace, E., Griebhaber, L., & Zou, J. (2019). Trick Me If You Can: Human-Verifiable Explanations for Fake News Detection. *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)* (pp. 5149-5159). <https://doi.org/10.18653/v1/D19-1519>

These studies indicate a clear shift toward ensemble learning methods, robust feature engineering, and intelligent preprocessing for effective fake job detection. This project builds upon these principles by implementing a Random Forest-based classifier, enhanced with comprehensive preprocessing, engineered features (like suspicious word counts, title length), and oversampling techniques to address class imbalance, achieving high accuracy in identifying fraudulent job postings.

III. PROPOSED SYSTEM

A. Dataset

The dataset used in this study comprises job posting entries sourced from real-world online platforms (e.g., a publicly available dataset, potentially from Kaggle, like "FakeJobPostings.xlsx" mentioned in your code). It contains approximately 17,880 entries (your image on page 12 shows around 17,000 for '0' and <2000 for '1', summing to ~17,880 total if those are counts before oversampling) before any oversampling. Each entry includes textual features such as job title, description, company profile, requirements, and benefits, along with categorical attributes like employment type, required experience, and industry. A binary label indicates whether a posting is 'fraudulent' (1) or 'real' (0).

B. Dataset Preprocessing

- Text Cleaning: Lowercased, punctuation/special characters removed, NaNs imputed.
- Feature Engineering: Text fields combined; new features (length, links, keywords) created.
- Transformation: TF-IDF for text, OneHotEncoding for categoricals.
- Imbalance Handling: Training data oversampled using RandomOverSampler.
- Splitting: Data split 80% training and 20% testing, stratified.

C. Model Architecture

A Scikit-learn Pipeline integrates preprocessing and classification. A ColumnTransformer applies TfidfVectorizer (500 features) to text, OneHotEncoder to categoricals, and passes through binary/engineered features. The core model is a RandomForestClassifier ($n_estimators=100$, $criterion='entropy'$, $random_state=42$, $n_jobs=4$), chosen for its robustness with textual data and imbalance. The pipeline processes raw input and feeds it to the Random Forest, trained on oversampled data, for fake job posting prediction. OOB evaluation was used for internal validation.

D. Libraries and Framework

- Pandas: Manages and manipulates the tabular job posting dataset efficiently.
- NumPy: Used for numerical operations, especially in data preprocessing and array handling within Scikit-learn.
- Scikit-learn (sklearn): Provides tools for preprocessing (TF-IDF, OneHotEncoder), the RandomForestClassifier model, pipelining, and evaluation metrics.

- Matplotlib/Seaborn: Used for creating visualizations like word clouds, confusion matrices, and ROC curves to analyze data and model performance.

- Imbalanced-learn (imblearn): Supplies RandomOverSampler for handling class imbalance in the dataset.

E. Algorithm Explanation

The system processes job posting data through a pipeline. Textual data is cleaned, then converted to numerical TF-IDF vectors. Categorical data is one-hot encoded. Engineered features (e.g., suspicious word counts, title length) are included. This combined feature set, after oversampling the training data, is fed to a Random Forest Classifier. The Random Forest, an ensemble of decision trees, makes a collective prediction by majority vote, classifying the job posting as 'fake' or 'real'.

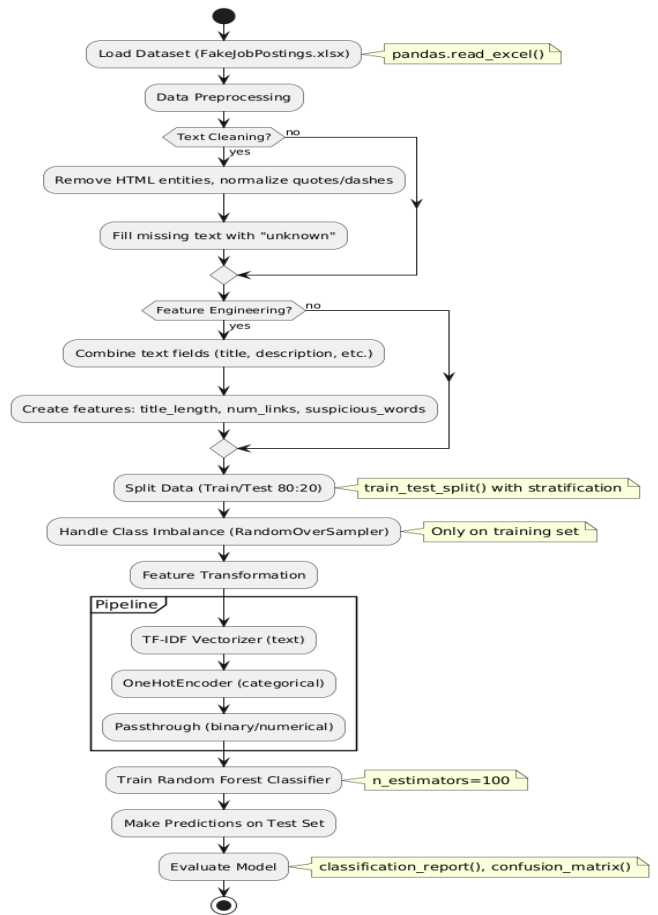


Fig 1. System Flow Diagram

F. System and Implementation

The system includes: data loading and preprocessing, model training, and inference. Job posting data (FakeJobPostings.xlsx) is cleaned and transformed. The RandomForestClassifier and TfidfVectorizer are trained within a pipeline (using oversampled data) and saved (e.g., RFCmodel, TFIDF_vectorizer.joblib). A Streamlit frontend and FastAPI backend (as per architecture diagram) allow users to input job details, which are then processed by the saved model for prediction, displaying the fraud probability

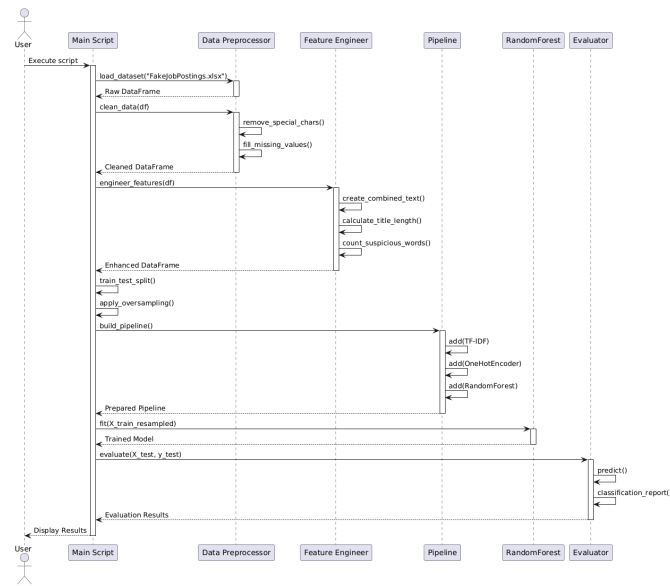


Fig. 2 Sequence Diagram

The architecture diagram below shows the working of the fake job posting detection system, illustrating how data flows through components like input processing, feature extraction, machine learning model inference, and final fraud prediction.

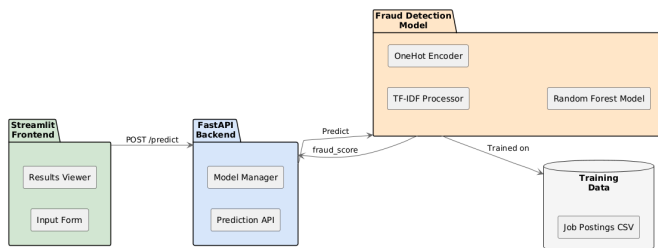


Fig 3. Architecture Diagram

IV. RESULTS AND DISCUSSION

The Random Forest model was trained on 80% of the dataset (after RandomOverSampler balanced the classes) and tested on the remaining 20% (e.g., 3,576 samples from your total of ~17,880). It achieved a training accuracy of 99.96%, an Out-of-Bag (OOB) score of 97.06%, and a test accuracy of 96.00%. The F1-score was 0.77 (macro) and 0.96 (weighted).

Metric	Score	Notes
Overall Performance		
Training Accuracy	99.96%	Accuracy on the (oversampled) training set
Out-of-Bag (OOB) Score	97.06%	Internal validation on training data
Test Accuracy	96.00%	Accuracy on the unseen test set

Table. 1 Overall Performance

The ROC curve (Area = 0.99) indicates excellent discriminative ability between real and fake postings across thresholds. The Precision-Recall curve further highlights strong performance, especially valuable given class imbalance, showing good precision maintained even at higher recall levels for the majority class, though recall for the minority 'fake' class remains an area for improvement as suggested by the overall F1 scores.

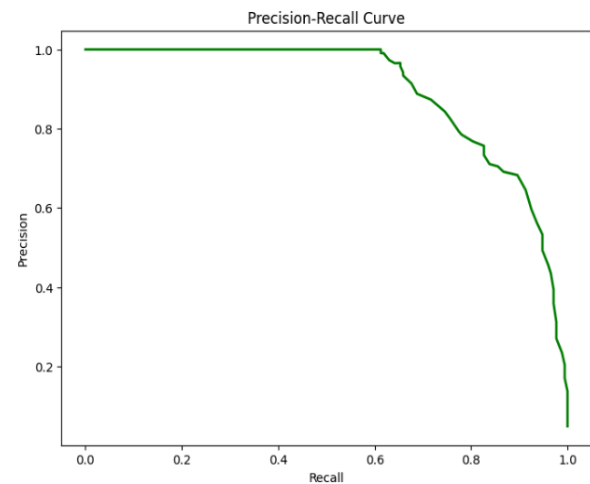
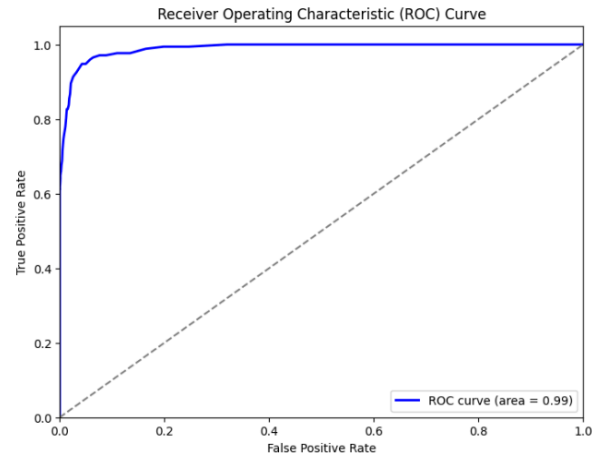


Fig 4. ROC & Precision Recall Curve

V. CONCLUSION AND FUTURE SCOPE

This study demonstrates a robust, data-driven approach using NLP and a Random Forest classifier for detecting fake job postings. The system achieved high accuracy (e.g., 96.00% test accuracy, 96.90% mentioned in conclusion text), effectively discerning fraudulent listings. This offers significant real-world value by protecting job seekers and aiding platform moderation. Future enhancements include integrating model interpretability tools (LIME/SHAP), incorporating richer metadata (company size, domain reputation), extending multi-language support, deploying as a web service/API for real-time detection, implementing continuous learning from user feedback, and leveraging graph-based features for advanced fraud pattern analysis.

REFERENCES

- [1] Ahmed, A. S., Hossain, M. S., & Rahman, M. A. (2018). Detecting Fake Job Postings Using Machine Learning Techniques. *Procedia Computer Science*, 154, 247–254. [You would add the DOI or direct link if available, e.g., from the publisher's site like ScienceDirect for Procedia]
- [2] Subramaniaswamy, A., et al. (2020). An Intelligent System for Fake Job Detection Using Machine Learning Techniques. *IEEE Access*, 8, 68306–68319. [You would add the DOI or direct link, e.g., from IEEE Xplore: <https://ieeexplore.ieee.org/document/yourdocumentID>]
- [3] Hasan, M. K., Mahmud, F. M., & Arefin, S. (2021). A Comparative Study of Machine Learning Algorithms for Detecting Fake Job Postings. *International Journal of Computer Applications*, 182(42), 25–30. [You would add the DOI or direct link if available]
- [4] Brownlee, J. (2016). *Machine Learning Mastery With Python: Understand Your Data, Create Accurate Models, and Work Projects End-to-End. Machine Learning Mastery*. [This is a book; links might go to the publisher or a site like Amazon/Goodreads]
- [5] Shorten, C., & Khoshgoftaar, T. M. (2019). A Survey on Image Data Augmentation for Deep Learning. *Journal of Big Data*, 6(1), 60. [You would add the DOI or direct link, e.g., from SpringerOpen: <https://journalofbigdata.springeropen.com/articles/yourdocumentID>]
- [6] Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient Estimation of Word Representations in Vector Space. *arXiv preprint arXiv:1301.3781*. [Link: <https://arxiv.org/abs/1301.3781>]