

Fake Job Posting Detection Using Machine Learning and Deep Learning

Title and Short Description

This project addresses the critical problem of detecting fake job postings, which have become increasingly prevalent in online job markets. Fake job ads not only waste job seekers' time but can also lead to financial scams and security risks. The importance of this problem lies in protecting vulnerable individuals from fraudulent schemes that exploit their job search efforts. Our approach uses Natural Language Processing (NLP) techniques combined with machine learning and deep learning models to automatically classify job postings as real or fake based on textual content.

We implemented and compared four different models - Naive Bayes, Logistic Regression, Support Vector Machine (SVM), and Artificial Neural Network (ANN) - on two distinct datasets. The results show that the ANN model achieved the highest accuracy (up to 98.4%) and AUC scores, demonstrating the effectiveness of deep learning approaches for this text classification task. The project provides both a comparative analysis of model performance and a practical web application for real-time job posting verification.

What Was Done in This Project

1. Data Collection and Preparation
 - Downloaded and analyzed two datasets: Fake Job Posting dataset from Kaggle and custom Job Train dataset
 - Performed extensive data cleaning, preprocessing, and feature engineering
 - Created comprehensive data pipelines for text processing
2. NLP Implementation
 - Built complete text preprocessing pipeline including tokenization, lemmatization, stopword removal
 - Implemented TF-IDF vectorization for feature extraction
 - Handled text normalization and cleaning for optimal model input
3. Model Development
 - Implemented four different machine learning models from scratch
 - Trained and fine-tuned each model with cross-validation
 - Saved trained models for deployment and future use
4. Performance Evaluation
 - Created detailed evaluation metrics and comparison tables
 - Generated visualizations including ROC curves, confusion matrices, and performance charts
 - Analyzed model strengths, weaknesses, and computational efficiency
5. Web Application Development

- Built a complete Flask-based web application with modern UI
- Implemented real-time prediction functionality
- Created interactive tabs for analysis and model insights
- Added responsive design with dark mode and animations

6. Comparative Analysis

- Performed thorough comparison across multiple datasets
- Analyzed model generalizability and robustness
- Documented findings with detailed explanations

7. Documentation and Deployment

- Created comprehensive README with all required sections
- Provided step-by-step installation and usage instructions
- Ensured code quality and reproducibility

Dataset Source

Primary Dataset: Fake Job Posting Dataset

Source: Kaggle (Real or Fake Job Posting Prediction dataset)

URL: <https://www.kaggle.com/shivamb/real-or-fake-fake-jobposting-prediction>

Size: Approximately 18,000 job postings

Features: Job title, location, department, company profile, description, requirements, benefits, telecommuting, has company logo, has questions, employment type, required experience, required education, industry, function, fraudulent (target variable)

Target Distribution: Binary classification (0 = Real, 1 = Fake)

Secondary Dataset: Job Train Dataset

Source: Custom compiled dataset from various job portals

Size: Approximately 10,000 job postings

Features: Similar to primary dataset but with different distribution

Purpose: Used to test model generalizability across different data sources

Data Preprocessing

- Data Cleaning: Removed duplicate entries and handled missing values
- Text Combination: Merged relevant text fields (title, description, requirements, company profile) into a single text feature
- Text Normalization: Converted to lowercase, removed special characters and punctuation
- Stopword Removal: Eliminated common English stopwords using NLTK
- Tokenization and Lemmatization: Broke text into tokens and reduced words to their base forms
- Feature Extraction: Applied TF-IDF (Term Frequency-Inverse Document Frequency) vectorization
- Final Dataset Size: 15,000–20,000 samples with TF-IDF features

Methods

Approach Overview

1. Data Acquisition and Preprocessing

2. Model Selection
3. Training and Validation
4. Performance Evaluation
5. Comparison and Analysis

Why This Approach?

- NLP techniques ensure efficient text representation
- Multiple models provide comparative insights
- Comprehensive metrics beyond accuracy
- Real-world deployable solution

Alternative Approaches Considered

- Traditional ML only (SVM, Random Forest)
- Advanced NLP (BERT, Transformers)
- Manual vs Automated Feature Engineering
- Ensemble Methods (not implemented)

Model Architectures

Naive Bayes

- Multinomial classifier assuming feature independence
- Fast and efficient baseline

Logistic Regression

- Linear model with sigmoid activation
- Regularization for overfitting control
- Interpretable and efficient

Support Vector Machine (SVM)

- Linear kernel optimized for text
- Excellent generalization on sparse data

Artificial Neural Network (ANN)

- Input: TF-IDF features
- Hidden Layers: 128 neurons, ReLU activation
- Dropout: 0.3
- Output: Sigmoid neuron
- Optimizer: Adam
- Loss: Binary cross-entropy
- Metrics: Accuracy, AUC

Steps to Run the Code

Prerequisites

Python 3.8+, Jupyter Notebook, and libraries: pandas, numpy, scikit-learn, tensorflow, keras, nltk, matplotlib, seaborn

Installation

```
git clone https://github.com/yourusername/fake-job-detection.git
```

```
cd fake-job-detection
```

```
pip install pandas numpy scikit-learn tensorflow keras nltk matplotlib seaborn
```

Execution Steps

1. Download datasets (fakejobposting.csv, job_train.csv)
2. Run notebooks:
 - o job_train.ipynb
 - o real_fake_job_posting.ipynb
3. Run web app:
 - o python app.py
 - o Access at <http://localhost:5000>

Pre-trained models are saved as .pkl and .h5 files

Experiments/Results Summary

Experimental Setup

Train-Test Split: 80-20

Cross-Validation: 5-fold

Hyperparameters:

- TF-IDF max_features=5000
- SVM C=1.0
- ANN epochs=50, batch_size=32
-

Performance Comparison

Job Train Dataset Results

| Model | Accuracy | AUC |
|---------------------|---------------|--------|
| ANN (Keras) | 97.31% | 0.9359 |
| SVM (Linear) | 96.92% | 0.8749 |
| Logistic Regression | 96.03% | 0.9441 |
| Naive Bayes | 90.49% | 0.9329 |

Fake Job Posting Dataset Results

| Model | Accuracy | AUC |
|---------------------|---------------|--------|
| ANN (Keras) | 98.40% | 0.9877 |
| SVM (Linear) | 98.29% | 0.9751 |
| Logistic Regression | 97.48% | 0.9839 |
| Naive Bayes | 91.55% | 0.9747 |

Visualization and Analysis

ANN achieved top accuracy and AUC.

ROC curves confirmed superior discriminative ability.

Confusion matrix showed ANN's balance between precision and recall.

SVM and Logistic Regression offered good alternatives for interpretability and efficiency.

Comparison with Published Methods

- ANN model outperforms prior ML-based studies

- Comparable to transformer-based models
- Multi-dataset evaluation ensures robust results

Hyperparameter Sensitivity

- ANN sensitive to learning rate, dropout
- SVM optimal with linear kernel
- Increasing TF-IDF features beyond 5000 yields minimal gain

Conclusion

1. ANN achieves superior performance (up to 98.4%)
2. Preprocessing significantly improves model results
3. Models generalize across datasets
4. SVM provides interpretable alternative
5. Flask web app adds practical usability

This project demonstrates AI's potential in combating online job fraud using NLP and Deep Learning. Future scope includes exploring transformers and larger datasets.

References

1. Kaggle Dataset: Real or Fake Job Posting Prediction
2. Pedregosa et al. (2011). Scikit-learn: Machine Learning in Python
3. Chollet, F. (2015). Keras
4. Bird et al. (2009). Natural Language Processing with Python
5. Zhang & Wallace (2015). CNN for Sentence Classification
6. Joachims, T. (1998). Text categorization with SVMs