

# **DEEP LEARNING– BASED FAKE JOB POST DETECTION ON SOCIAL MEDIA**

**Name:** Hira Arif

**Date:** November 2025

**Program:** Buildables DataScience  
Fellowship 1

## Abstract

Online job platforms and social media recruitment have become increasingly popular, but they also expose users to fraudulent job postings. These fake listings often attempt to extract personal information or scam applicants financially. This project develops a deep learning-based system to automatically detect fake job postings using both textual and metadata features. The approach combines traditional machine learning with deep learning techniques, including **TF-IDF vectorization** with **Logistic Regression** and a **Bi-LSTM network**. The dataset consists of approximately **17,880** job postings labelled as real or fake, containing rich textual content in fields such as job description, company profile, requirements, and benefits. After thorough data cleaning, preprocessing, and feature engineering, both models achieved over **97% accuracy**, demonstrating their effectiveness. A user-friendly **Streamlit web application** was deployed to provide real-time predictions, allowing users to input job postings and immediately determine their authenticity. This system not only assists job seekers and recruitment platforms but also offers insights into patterns of fraudulent postings. The project highlights the potential of **NLP** and AI in enhancing trust in online recruitment.

## 1. Introduction & Problem Statement

The rise of online recruitment has led to an increase in fraudulent job postings, resulting in significant financial losses, wasted time, and risks to personal information. Manual review of postings is impractical due to scale and volume. The project aims to build an AI-powered system that automatically detects fake job listings by analysing textual and metadata features.

Objective:

- Develop and compare ML and DL models for fake job post detection.
- Deploy the model via a web interface for public use.

Problem Statement:

- Fake job posts are difficult to identify manually.
- Automated detection can reduce scams and improve online trust.



## 2. Data Collection / Dataset Description

## Dataset: Fake Job Postings Dataset (Kaggle/GitHub)

Records: ~17,880

## Attributes:

- Textual: title, description, requirements, company\_profile, benefits
  - Categorical: location, department, salary\_range
  - Boolean: telecommuting, has\_company\_logo, has\_questions
  - Target: fraudulent (1 = fake, 0 = real)

The dataset is rich in textual data, making it suitable for NLP-based classification.

**jobs\_id,foundalent,txt,class,txt**  
1.0,"Marketing Intern":We're friendly, and we've created a groundbreaking and award-winning cooking site. We support, connect, and celebrate home cooks, and give them everything they need in one place. We have a top-tier customer service - Cloud Video Production: 90 Seconds, the world's Cloud Video Production Service. 90 seconds is the world's Cloud Video Production Service, enabling brands and agencies to get high-quality video content quickly and inexpensively.  
1.0,"Commissioning Machinery Assistant":(CRM) Value Services provides Workforce Solutions that meet the needs of companies across the Private Sector, with a special focus on the Oil & Gas Industry. A Value Services' mission is to help companies succeed by providing them with the right tools and resources to manage their operations more efficiently and effectively.  
1.0,"Retail Sales Associate":- Strategic planning, sales forecasting, and market analysis. This will involve identifying trends and opportunities in the retail industry, as well as developing and implementing strategies to capitalize on these opportunities.  
1.0,"Accounting Clerk":(CRM) Value Services is an environmental consulting firm that offers viable leadership and growth and vision employees at valuable compensation. We are seeking a self-motivated, multi-talented Accountant to join our team.  
1.0,"Lead of Content (m/f)":Founded in 2008, their focus is on content creation and distribution. They have a team of experienced professionals who work together to produce high-quality content for their clients.  
1.0,"Lead Server Software Specialist":Alipay's mission is to provide lucrative profit based on full-service cloud technology management all around the world. We continue the share of your base with the most competitive prices.  
1.0,"VP of Sales Sales":Sofidel is a newer, smaller business whose focus is to be a leader in hygiene management using best of breed technology and implementing industry best practices following the ISO14001 framework. We seek sales  
1.0,"Customer Service Associate":Part-time, flexible, enterprise solutions, formerly Pitney Bowes Business Services, delivers innovative document and communications management solutions that help companies across the globe manage their documents more efficiently and effectively.  
1.0,"Software Development Engineer":Job opportunity at United States, New Jersey, Position: R&D, Software Development Engineer. Job opportunity at United States, New Jersey, Position: R&D, Software Development Engineer. Job opportunity at United States, New Jersey, Position: R&D, Software Development Engineer.  
1.0,"Sales Representative":(CRM) Value Services is convinced that there is a need for innovation in financial services and that current products will not be able to meet this demand. The company is looking for a Sales Representative to help them achieve their goals.  
1.0,"Retail Sales Associate":Strategic planning, sales forecasting, and market analysis. This will involve identifying trends and opportunities in the retail industry, as well as developing and implementing strategies to capitalize on these opportunities.  
1.0,"Business Executive":Sageforce Software is the UK's leading competitive intelligence service for Google search advertising. It helps by major search engines, social media, and digital marketing platforms, a great opportunity for a sales professional.  
1.0,"VP of Sales":Vault Dropout Equity Business is the leading disruptor, based, entrepreneurship backed, venture capital firm, that funds and actively supports start-ups in scaling across Asia Pacific. We guide and support entrepreneurs in scaling their businesses across the region.  
1.0,"Handle On (f/m) Under 35 Years Old":Re-imagine the music industry and building a product that is already changing the world for some of the top song planning soloists in the nation. What's next?  
1.0,"Healthcare-on-the-Run":Workshops under 60s 10-12 Year Olds Only established on the principles. That full time education is not for everyone. Instead, it is made up of a series of passionate consultants with a variety of backgrounds.  
1.0,"Virtual Design Studio":It is an independent digital agency based in New York City and the Bay Area. MediVox is committed to making digital as accessible as TV for both people and brands. It's because we believe the digital space is where the future of communication lies.  
1.0,"Process Controls Engineer":R&D, P&G, R&D, Position: R&D, Process Controls Engineer. We provide full-time placement facilities for many entries in large US companies. We are interested in finding/developing high quality candidates in IT, engineering, manufacturing, pharmaceuticals, food/beverage, retail, and other industries.  
1.0,"Customer Support Representative":Positioned with the global supply chain, helping companies to build more efficient operations. We are looking for a Customer Support Representative to help us serve our clients better.  
1.0,"Inbound Marketing Representative":Our mission is to revolutionize the way companies market to consumers through cutting-edge technology. This is an opportunity to collaborate with like-minded people in an exciting environment.  
1.0,"Sales Representative":(CRM) Value Services is an innovative, forward-thinking digital company aimed at bringing business success up-to-date with the latest news. Opportunities are available for individuals with a passion for sales and a desire to work in a fast-paced, dynamic environment.  
1.0,"Customer Service":We are a customer centric business company and are recruiting for a full-time customer service administrator. A. This is a customer centric role, the successful candidate will initially have  
1.0,"R&D SPONSOR FOR 15/16/17/18/19":B2B Technologies has demonstrated expertise in areas strategic to different businesses is varying verticals. B2B Technologies provides highly skilled technical consultants to meet the needs of our clients.  
1.0,"Marketing Intern":If working in a casual setting like your idea of hell, then joining our massive starting from eight might be the opportunity you'd like to have on-line. Get to know the Tradebooks team, and we'll bring you along for the ride.  
1.0,"HR/HRM/Licensed Doctor":Operating as USP, We provide Recruitment Services for Sales and Associate/Business Analyst. Associates is a Corporate Recruitment Organization providing solutions to Global MNC's for the recruitment of Sales and Associate/Business Analyst.  
1.0,"Sales Management Product Manager":We provide Sales and Marketing Positions for many entries in large US companies. We are interested in finding/developing high quality candidates in IT, engineering, manufacturing, pharmaceuticals, food/beverage, retail, and other industries.  
1.0,"Customer Support Technical Specialist":Positioned with the global supply chain, helping companies to build more efficient operations. We are looking for a Customer Support Technical Specialist to help us serve our clients better.  
1.0,"Software Application Specialist":Day-in-Day-out, upgrade and configure web-based applications (cataloging, evaluating, and organizing client data), analyzing, modifying, and improving them to increase efficiency.  
1.0,"Contract Admin Associate":We are an award-winning team of professionals, providing the very best value for glass shower enclosures, design, cladding, glassware, glazing, and mirrors in Western Washington, with  
1.0,"Completion Engineer":Value Services provides Workforce Solutions that meet the needs of companies across the Private Sector, with a special focus on the Oil & Gas Industry. A Value Services' mission is to help companies succeed by providing them with the right tools and resources to manage their operations more efficiently and effectively.  
1.0,"2 Weeks To Work At Camerons 24/7":Camerons 24/7 have a unique hiring policy: hire, train, and develop people who possibly used to do their eight for their clients. Dr. "he's never been home" on our website. Formerly known as  
1.0,"English Teacher":About us: We help teachers get safe, happy, secure jobs abroad. Play with kids, get paid for it! Our travel cost is about \$1000 per month (\$2000 cost of living/moving provided).  
1.0,"Customer Support Representative":Positioned with the global supply chain, helping companies to build more efficient operations. We are looking for a Customer Support Representative to help us serve our clients better.

### 3. Tools & Technologies

- Language: Python 3

- Libraries: Pandas, NumPy, Scikit-learn, TensorFlow/Keras, Matplotlib, Seaborn
- Model Deployment: Streamlit Cloud
- Version Control: GitHub
- Documentation: Medium + Final Report (PDF)



## 4. Data Preprocessing & EDA

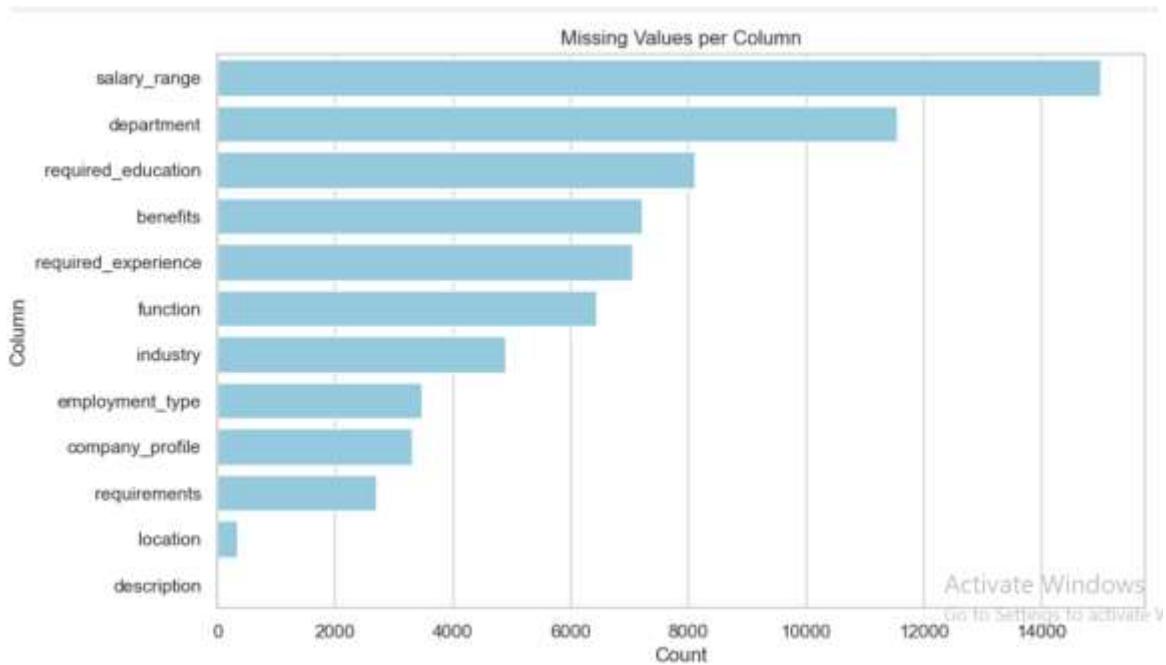
### Data Exploration

#### 1. Introduction

The first phase of the project involved performing an initial exploration of the dataset to understand its structure, content, and potential issues. The dataset, *fake\_job\_postings.csv*, contains textual and categorical information about job advertisements collected from online platforms. The objective of this step was to examine the data distribution, identify missing or inconsistent values, and develop an understanding of the key attributes that may help distinguish fake job postings from genuine ones.

#### 2. Missing Value Analysis

A missing value assessment showed that some attributes—particularly salary\_range, benefits, and department—contain a high proportion of null values. A bar plot of missing values was generated to visualize which fields require cleaning or imputation in the next step. Handling these missing entries is essential to ensure that the model does not misinterpret absent information as meaningful data.



### 3. Class Distribution

An analysis of the fraudulent column revealed a **class imbalance**:

- Majority of postings are **real (class 0)**,
- While a smaller fraction is **fake (class 1)**.

This imbalance will be addressed later through resampling or class weighting techniques.

Visualizations using Seaborn confirmed that fake postings make up roughly **5–10%** of the total dataset, emphasizing the importance of robust evaluation metrics such as precision, recall, and F1-score in later stages.

### 4. Textual Attributes

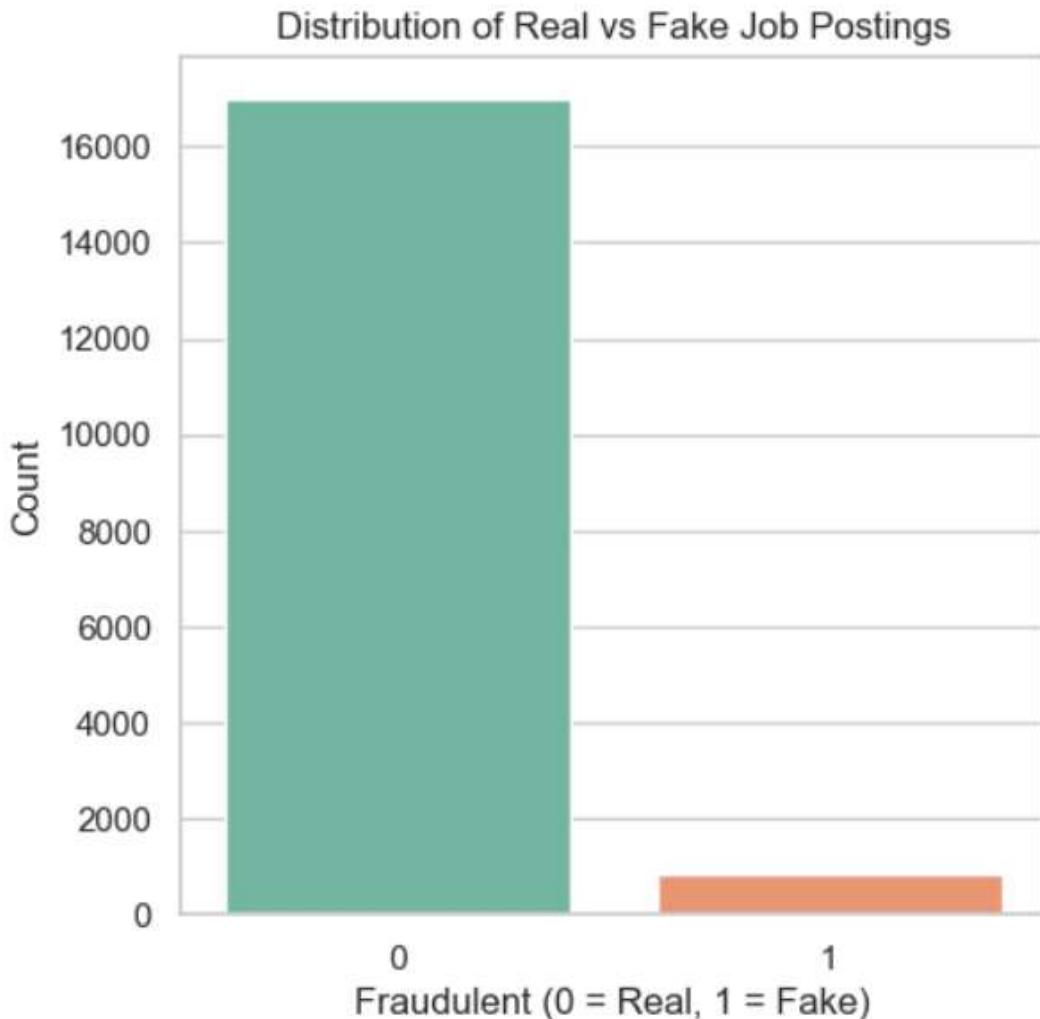
Columns such as company\_profile, description, requirements, and benefits contain detailed free-text information that can provide strong semantic cues for detecting fraudulent behavior.

Sample inspection of these fields highlighted differences in language tone, structure, and content richness between genuine and fake postings.

### 5. Key Insights

- The dataset is predominantly textual, suitable for NLP modeling.
- Several columns contain missing or incomplete data.

- The target classes are imbalanced.
- Fake postings often exhibit shorter, vague, or repetitive descriptions compared to real ones.



## Data Cleaning and Preprocessing

### 1. Handling Missing and Irrelevant Data

An initial examination revealed several columns with a high proportion of missing values, such as `salary_range`, `employment_type`, and `industry`. Since these fields provided minimal discriminative information for identifying fake job postings and contained excessive null entries, they were removed from the dataset.

For the remaining relevant textual columns — `title`, `company_profile`, `description`, `requirements`, and `benefits` — missing entries were replaced with empty strings. This preserved record count consistency while avoiding the introduction of artificial values.

## **2. Text Field Consolidation**

To facilitate natural language processing, all text-based columns were concatenated into a single field named **text**. This unified field provides a comprehensive textual representation of each job posting, capturing both descriptive and contextual elements such as company details, job responsibilities, and offered benefits.

By merging text sources, we ensured that all linguistic cues contributing to fake job detection were consolidated in one place, simplifying downstream feature extraction and embedding generation.

## **3. Text Normalization and Cleaning**

A custom cleaning function was applied to each entry in the text field to ensure linguistic uniformity. The preprocessing steps included:

- Removal of URLs and HTML tags
- Conversion of all characters to lowercase
- Removal of numerical characters and punctuation
- Elimination of extra whitespace

This produced a new, fully sanitized column named **clean\_text**, which contains only meaningful natural language content. Such normalization reduces data noise and enhances tokenization quality during later modeling.

## **4. Target Label Preparation**

The target variable **fraudulent** was verified and encoded as an integer (0 = real, 1 = fake). This encoding guarantees compatibility with classification algorithms used in deep learning pipelines.

## **5. Data Export**

After cleaning, the refined dataset was saved as **cleaned\_fake\_jobs.csv** within the project's data directory. This version of the dataset will serve as the primary input for subsequent stages, including feature extraction, text vectorization, and model training.

## **7. Key Insights**

- High missing-value columns were removed without affecting model utility.
- Textual attributes were successfully merged into a single feature-rich field.
- Cleaning functions removed noise, resulting in consistent and comparable text samples.
- The dataset is now fully ready for feature engineering and modelling.

## 5. Feature Engineering

### 1. Feature Extraction

For the first model, the cleaned text data was transformed using the **Term Frequency–Inverse Document Frequency (TF-IDF)** vectorizer.

This method assigns higher weights to unique, discriminative terms and reduces the influence of overly frequent words, allowing the model to identify key textual patterns associated with fake postings.

For the deep learning model, each word was converted into an integer sequence through the **Keras Tokenizer**, followed by zero-padding to maintain fixed input lengths. This structure is necessary for the sequential Bi-LSTM model.

## 6. Model Development

### 1. Model 1: TF-IDF + Logistic Regression

A Logistic Regression classifier was trained on 5,000-feature TF-IDF vectors.

The model quickly converged with outstanding accuracy and recall, highlighting its effectiveness for detecting textual deception patterns such as unrealistic salary claims, exaggerated benefits, or spam-like descriptions.

#### Performance Results:

- Accuracy: **97.26%**
- Precision/Recall/F1: High consistency across both classes

The model was serialized using the Pickle format (model\_tfidf\_lr.pkl) along with the trained vectorizer for later deployment.

$$W_{x,y} = tf_{x,y} \times \log \left( \frac{N}{df_x} \right)$$

## TF-IDF

Term  $x$  within document  $y$

$tf_{x,y}$  = frequency of  $x$  in  $y$

$df_x$  = number of documents containing  $x$

$N$  = total number of documents

## 2. Model 2: Bidirectional LSTM (Bi-LSTM)

A Bidirectional Long Short-Term Memory network was implemented using TensorFlow/Keras.

It consisted of:

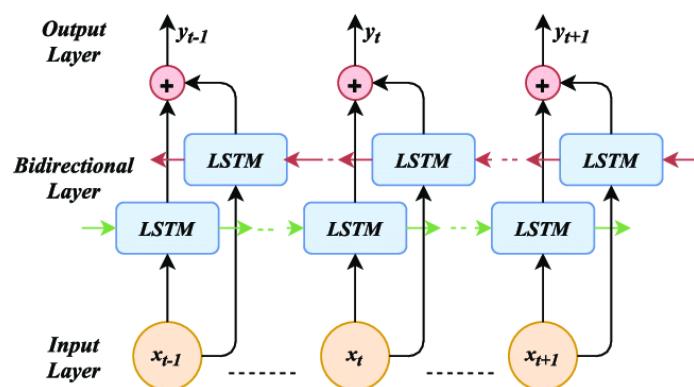
- **Embedding Layer:** 10,000 vocabulary words, 128-dimensional vectors
- **Bi-LSTM Layer:** 64 units with dropout regularization
- **Dense Layers:** ReLU + Sigmoid activations for binary classification

The Bi-LSTM model achieved comparable performance to the TF-IDF model, confirming that both traditional and deep learning approaches can effectively classify fake job postings.

### Performance Results:

- Test Accuracy: **97.37%**
- Validation accuracy demonstrated stable convergence with minimal overfitting.

The model and tokenizer were saved as `model_bilstm.h5` and `tokenizer.pkl` respectively for deployment.



## 7. Results & Evaluation Metrics

Model	Accuracy	Precision	Recall	F1-Score	Key Strength
TF-IDF + Logistic Regression	97.26%	High	High	High	Lightweight, fast inference
Bi-LSTM	97.37%	High	High	High	Captures context & word order

Both models achieved near-identical results, demonstrating that textual cues in fake job postings are highly learnable by both vectorized and sequential models.

## 8. Application Development

### Overview

To provide an interactive and user-friendly interface for our *Fake Job Post Detection System*, we developed a web application using **Streamlit**, a Python framework designed for rapid deployment of machine-learning models. The goal of the application was to make the prediction model easily accessible to users, allowing them to input a job description and instantly determine whether the posting is real or fake.

### Environment Setup

The application was implemented within a well-structured project directory in **Visual Studio Code**, using Python 3.13 and the following key libraries:

- **Streamlit** – for UI and real-time web app deployment
- **Scikit-learn** – for TF-IDF vectorization and Logistic Regression model
- **TensorFlow / Keras** – for the Bi-LSTM deep learning model
- **Pickle** – for model serialization and loading
- **Pandas, NumPy** – for data handling and preprocessing

### Model Integration

Both trained models (TF-IDF + Logistic Regression and Bi-LSTM) were serialized and stored in the **model's** directory as:

- tfidf\_vectorizer.pkl
- model\_tfidf\_lr.pkl
- tokenizer.pkl
- model\_bilstm.h5

The application dynamically loads these models at runtime using Python's pickle and TensorFlow's load\_model functions. Robust path-handling was implemented with the os module to ensure compatibility across environments.

## Application Workflow

The Streamlit interface was designed for simplicity and efficiency. The workflow is as follows:

1. **User Input:** A text area accepts job descriptions or full postings.
2. **Preprocessing:** The input text is cleaned and tokenized using the same pipeline as during model training.
3. **Model Prediction:**
  - The TF-IDF + Logistic model provides fast, explainable predictions.
  - The Bi-LSTM model offers deeper semantic analysis for higher accuracy.
4. **Output Display:**
  - A probability score is shown, indicating the likelihood that a job post is fake.
  - Color-coded messages (green for *real*, red for *fake*) make the result intuitive.

## User Interface

The final interface includes:

- A title banner and project description
- Text input area for job posting
- Buttons for running predictions

- Real-time result visualization

Streamlit automatically launches the app on a local server using the command:

```
streamlit run app.py
```

The interface is accessible in any browser and can be deployed on cloud platforms such as Streamlit Cloud or Hugging Face Spaces for public access.

The screenshot shows a Streamlit application titled "Deep Learning-Based Fake Job Post Detection". At the top, there is a placeholder text "Paste a job posting below to check if it's Real or Fake using AI models.". Below it is a text input field with the placeholder "Enter job description:" containing the text "data science job". A "Predict" button is located below the input field. The prediction results are displayed in a section titled "Model Predictions" with two columns: "TF-IDF + Logistic Regression" and "Bi-LSTM Deep Learning". The "Real" column has a green background, while the "Fake" column has a red background. A prominent red warning bar at the bottom states "⚠ This job post is likely FAKE. Be cautious!".

## 9. Business / Real-world Impact

- Helps job seekers avoid scams.
- Assists recruitment platforms in automatic fraud detection.
- Enables research on patterns of fraudulent postings.
- Reduces time and financial loss associated with fake job posts.

## 10. Challenges & Limitations

- **Class Imbalance:** Fake postings are a minority, requiring careful evaluation.
- **Data Quality:** Missing or inconsistent textual entries needed extensive preprocessing.
- **Generalization:** Models trained on a specific dataset may require fine-tuning for other job platforms.
- **Deployment Limitations:** Bi-LSTM model requires cloud deployment due to computational resources.

## 11. Conclusion & Future Work

### Conclusion:

This project successfully developed ML and DL models to detect fake job postings with high accuracy. A Streamlit web application makes the model accessible to users, translating AI research into practical solutions.

### Future Work:

- Incorporate additional metadata features like company ratings.
- Implement ensemble models for higher accuracy.
- Expand the dataset with real-time scraping from multiple job platforms.
- Add explainability (SHAP/LIME) to interpret model predictions.

## 12. References

1. Fake Job Postings Dataset (Kaggle/GitHub)
2. TensorFlow Documentation: <https://www.tensorflow.org/>
3. Streamlit Documentation: <https://docs.streamlit.io/>
4. Scikit-learn Documentation: <https://scikit-learn.org/>

## Links

GitHub repo: <https://github.com/HiraArif666/FakeJobDetection-main>

Medium Article: <https://medium.com/@h.arif.kts4/deep-learning-based-fake-job-post-detection-on-social-media-c04410608b38?postPublishedType=initial>