Executive Review

: [Spring 2024]

: [Linh Bui]

: [Fake Job Posting Detection]

## Summary:

My project aims to develop a classification model for identifying fraudulent job postings. The primary objectives include enhancing user trust, reducing the risk of scams, and maintaining the integrity of the job marketplace. By leveraging advanced machine learning algorithms, the project seeks to analyze textual and contextual features to accurately classify job postings as either legitimate or fake. Potential outcomes involve a significant decrease in fraudulent activities, improved user satisfaction, and increased platform credibility. Key metrics for success include precision, recall, and F1 score, with a target of achieving high accuracy in identifying fake job postings while minimizing false positives and negatives. This business goals of this project are providing a secure and trustworthy environment for both job seekers and employers, ultimately fostering a positive user experience, sustaining long-term platform growth, and finally, benefiting employment agencies by not only enhancing customer trust but also significantly reducing their workload and minimizing the risk of being scammed.

## Timeline:

Outline the project timeline with key milestones and deadlines. Include information on the current project status, any delays encountered, and adjustments made to the original schedule.

| Milestone | Planned Date |
|---|---|
| Project Kickoff | 01/30/2024 |
| Data Cleaning and Exploration | 02/08/2024 |
| Model Development | 03/07/2024 |
| Model Training and Optimization | 04/04/2024 |
| Final Model Deployment | 04/25/2024 |
| Performance Monitoring and Project Review | 05/02/2024 |
| Final Presentation | 05/08/2024 |

## Data Needs:

The project utilizes the "Real / Fake Job Posting Prediction" dataset from Kaggle, containing 18,000 job descriptions, with approximately 800 are fake. The dataset includes various features such as job title, location, salary range, company profile, description, requirements, benefits, and several others. The data types encompass both textual information and meta-information about the jobs, containing a mix of text strings, float, and integer types. While the advantage of not having to collect data from scratch alleviates concerns about data collection, challenges in data processing and cleaning still exist for the "Fake Job Posting Detection" project. One notable challenge lies in addressing missing values within the dataset, ensuring that crucial features are complete and representative. Additionally, since the dataset contains a mix of text strings, float, and integer values, harmonizing these diverse data types requires careful preprocessing. Imbalanced classes, particularly with a relatively small proportion of fake job postings, pose challenges in effectively handling and mitigating bias during model training. The quality of textual data, such as job descriptions, introduces complexities in natural language processing, demanding thorough text cleaning and feature extraction. These challenges collectively emphasize the importance

of meticulous data preprocessing to enhance the reliability and effectiveness of the machine learning model for fake job posting detection.

## Data Sources:

- [Source 1]: https://www.kaggle.com/datasets/shivamb/real-or-fake-fake-jobposting-prediction
- [Source 2]: I will collect more data from job posting sites if needed.

## Data Challenges:

- [Challenge 1]: The dataset's imbalance, with only around 800 fake job postings out of 18,000, poses a challenge in training a balanced model. The problem with a classification model trained on imbalanced data is that the model learns that it can achieve high accuracy by consistently predicting the majority class, even if recognizing the minority class is equal or more important.
- [Challenge 2]: Processing textual data in features like "description" demands careful handling of natural language complexities. Cleaning, stemming, and vectorization are essential for transforming text into a format suitable for machine learning. Coping with diverse writing styles, potential misspellings, and nuances requires meticulous preprocessing to ensure the model captures meaningful patterns indicative of fake job postings.

## Value Proposition:

My project delivers substantial business value by enhancing trust, reducing operational workload, and safeguarding the integrity of the job marketplace. Through the implementation of a classification model, we anticipate a significant decrease in fraudulent activities, directly impacting key performance indicators (KPIs) such as user satisfaction and platform credibility. Employment agencies stand to benefit from reduced verification efforts, leading to operational efficiency and cost savings. Additionally, the project aligns with the business goal of providing a secure environment, fostering increased user confidence, which can potentially drive user acquisition and retention. The reduction in the prevalence of fake job postings contributes to a positive user experience, supporting long-term platform growth and revenue generation. Overall, the value proposition lies in creating a safer, more trustworthy job marketplace, positively impacting KPIs, operational efficiency, and user engagement.

## Key Benefits:

- [Benefit 1]: Enhance user trust with a secure job marketplace. The project ensures authenticity, providing a reliable platform for both job seekers and employers.
- [Benefit 2]: Streamline operations for employment agencies in a safe job market. The model's automated detection of fake job postings improves efficiency, allowing agencies to save verification cost/efforts and focus on genuine opportunities and reducing risks associated with fraudulent listings.

## KPI Impact:

As this is only a simulated project addressing a real-world issue, it is challenging for me to accurately assess the KPI impact. For instance, to evaluate user trust, I would need to conduct surveys or gather feedback; to measure the operational efficiency, I would need to analyze the time and resources spent by employment agencies on job verification, but such real-world assessments are beyond the scope of this simulation. Thus, the only indicator that I can measure is the number of accurately detected fake job postings out of the total of 800.

| KPI | Before Project | After Project | Improvement |
|---|---|---|---|
| **The number of accurately detected fake job postings** | 0 | | |

## Elevator Pitch:

In a crowded job marketplace, my 'Fake Job Posting Detection' data science project is the solution to a pervasive issue – fraudulent job listings. By leveraging advanced machine learning, I've built a classification model that accurately identifies fake postings, creating a secure job environment. The result? Increased user trust, a safer job market, and streamlined operations for employment agencies. My solution not only safeguards the integrity of the job platform but also fosters confidence among job seekers and employers, driving user satisfaction and operational efficiency. It's not just about detecting fakes; it's about building a trusted, reliable, and efficient job marketplace for everyone.

| KPI | Before Project | After Project | Improvement |
|---|---|---|---|
| **The number of accurately detected fake job postings** | 0 | | |