

IE 5390 – Final Case – Fraud detection Case Study

Name: Pooja Arumugam

Final Project Report: Fake Job Postings Analysis

1. Introduction

The internet has revolutionized job searching, providing job seekers with thousands of opportunities at their fingertips. However, alongside genuine job postings, fraudulent job listings have emerged, exploiting unsuspecting candidates. Fake job postings can lead to identity theft, financial loss, and even personal safety risks.

This analysis aims to detect fraudulent job postings by identifying key differences between legitimate and fake job listings. Using data analytics, we can uncover hidden patterns that may indicate fraud and develop a methodology to flag suspicious postings.

A) Scenario Identification

- **Project Objective**

Fraudulent job postings are a growing concern as they can deceive job seekers, leading to identity theft, financial fraud, and wasted time. This project aims to analyze job posting data to detect patterns that distinguish fake job postings from legitimate ones.

- **Research Questions (Questions to answer):**
 1. *What factors differentiate fake job postings from real ones?*
 2. *Are there common attributes among fraudulent job listings?*
 3. *Can data analysis techniques effectively detect fake job postings?*
- **Significance of the Analysis (Why is it important)?**

This analysis is important because:

- It helps job seekers recognize fraudulent job postings and avoid scams.
- It provides insights for job boards and recruitment platforms to enhance their fraud detection systems.
- It explores how machine learning and statistical techniques can improve the detection of fraudulent job listings.

B) Dataset description:

Description of Variables or key elements

- The dataset consists of multiple columns such as job title, company, location, job description, required experience, employment type, fraud label (0 = Real, 1 = Fake), and other metadata related to job postings.
- This dataset contains 18K job descriptions out of which about 800 are fake. The data consists of both textual information and meta-information about the jobs. The dataset can be used to create classification models which can learn job descriptions which are fraudulent.

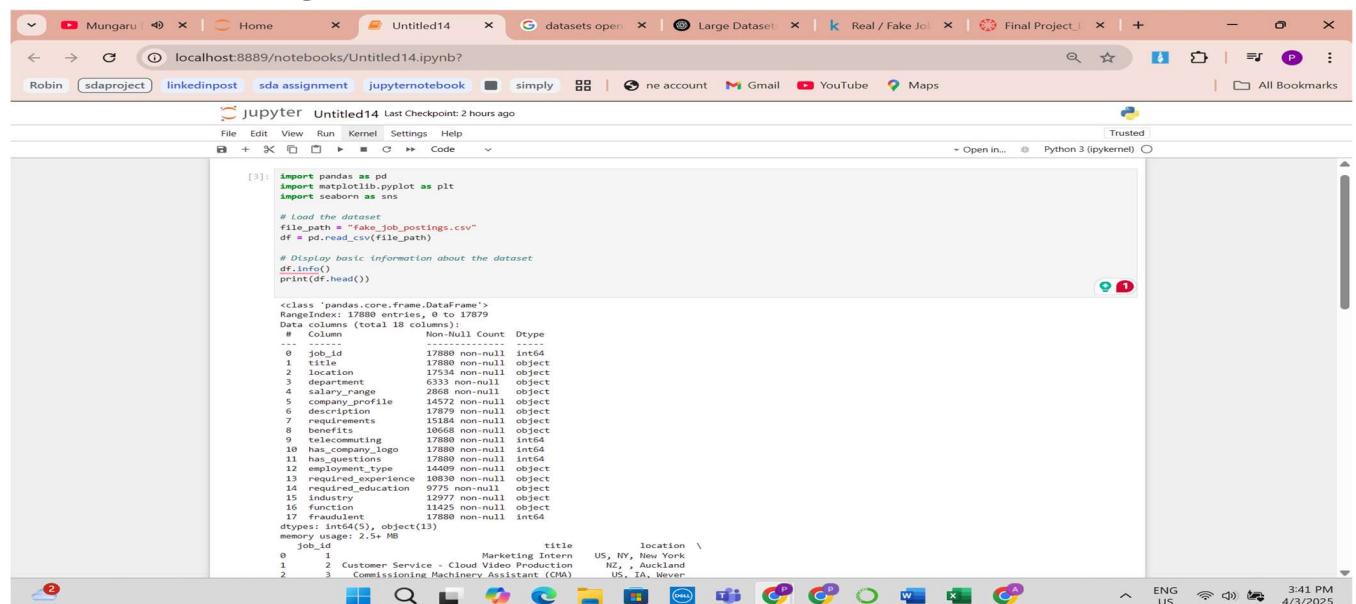
C) Dataset Source and Dataset File

- The dataset used for this analysis is fake_job_postings.csv.
- Source: <https://www.kaggle.com/datasets/shivamb/real-or-fake-fake-jobposting-prediction>

D) Characteristics of the Data (As-Is Form)

- The dataset contains 17,880 job postings with 18 columns related to job details, company information, and fraud labels.
- The fraudulent column is a binary variable (0 = Legitimate, 1 = Fraudulent).
- Key job attributes include title, location, company profile, job description, requirements, benefits, employment type, industry, and function.

Data cleaning:



A screenshot of a Jupyter Notebook interface running in a browser window. The notebook has a single cell containing Python code for loading a dataset and displaying its first few rows. The code is as follows:

```
[3]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# Load the dataset
file_path = "fake_job_postings.csv"
df = pd.read_csv(file_path)

# Display basic information about the dataset
df.info()
print(df.head())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 17880 entries, 0 to 17879
Data columns (total 18 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   job_id          17880 non-null   int64  
 1   title           17880 non-null   object 
 2   location        17534 non-null   object 
 3   department      6333 non-null   object 
 4   size            2686 non-null   object 
 5   company_profile 14572 non-null   object 
 6   description     17879 non-null   object 
 7   requirements    15184 non-null   object 
 8   benefits         10620 non-null   object 
 9   telecommuting   17880 non-null   int64  
 10  has_company_logo 17880 non-null   int64  
 11  has_full_time   17880 non-null   int64  
 12  employment_type 14409 non-null   object 
 13  required_experience 10830 non-null   object 
 14  required_education 10830 non-null   object 
 15  industry         12977 non-null   object 
 16  function         11425 non-null   object 
 17  salary           17880 non-null   int64  
dtypes: int64(5), object(13)
memory usage: 2.5+ MB
```

The output of the code shows the first few rows of the dataset:

job_id	title	location
1	Marketing Intern	US, NY, New York
2	Customer Service - Cloud Video Production	NZ, Auckland
3	Commissioning Machinery Assistant (CMA)	US, IA, Waver

Figure 1: Data Cleaning

```

jupyter Untitled14 Last Checkpoint: 2 hours ago
File Edit View Run Kernel Settings Help Trusted
+ % ▶ Code ▶ Open in... Python 3 (ipykernel)
 1/1889 non-null int64
dtypes: int64(5), object(13)
memory usage: 2.5+ MB
job_id      title          location \
0          1  Marketing Intern   US, NY, New York
1          2  Customer Service - Cloud Video Production  NZ, , Auckland
2          3  Commissioning Machinery Assistant (CMA)  US, IA, Waver
3          4  Account Executive - Washington DC  US, DC, Washington
4          5  Bill Review Manager  US, FL, Fort Worth

department salary_range           company_profile \
0  Marketing          NaN  We're Food52, and we've created a groundbreakin...
1  Success           NaN  99 Seconds, the world's Cloud Video Production ...
2  NaN               NaN  Valor Services provides Workforce Solutions th...
3  Sales             NaN  Our passion for improving quality of life thro...
4  NaN               NaN  SpotSource Solutions LLC is a Global Human Cap...

description \
0  Food52, a fast-growing, James Beard Award-winn...
1  Organised - Focused - Vibrant - Awesome! Our ...
2  Our client, located in Houston, is actively se...
3  THE COMPANY: ESRI - Environmental Systems Rese...
4  JOB TITLE: Itemization Review Manager LOCATION:...

requirements \
0  Experience with content management systems a m...
1  What we expect from you: Your key responsibilit...
2  Implement pre-commissioning and commissioning ...
3  EDUCATION: Bachelor's or Master's in GIS, busi...
4  QUALIFICATIONS: RN license in the State of Texa...

benefits telecommuting \
0          NaN  0
1  What you will get from usThrough being part of...  0
2          NaN  0
3  Our culture is anything but corporate-we have ...  0
4          Full Benefits Offered  0

has_company_logo has_questions employment_type required_experience \
0          1          0        Other    Internship
1          1          0      Full-time  Not Applicable
2          1          0          NaN        NaN
3          1          0      Full-time  Mid-Senior level
4          1          1      Full-time  Mid-Senior level

```

Figure 2

```

jupyter Untitled14 Last Checkpoint: 2 hours ago
File Edit View Run Kernel Settings Help Trusted
+ % ▶ Code ▶ Open in... Python 3 (ipykernel)
name      department      title
1  NaN  Marketing and Advertising  Marketing
2  NaN  Computer Software          Customer Service
3  Bachelor's Degree  Hospital & Health Care  Sales
4  Bachelor's Degree  Hospital & Health Care  Health Care Provider

fraudulent
0  0
1  0
2  0
3  0
4  0

[5]: # Check for missing values
missing_values = df.isnull().sum().sort_values(ascending=False)
missing_values = missing_values[missing_values > 0]

[7]: missing_values
salary_range      16912
departments       11547
required_education  8105
benefits          7212
required_experience  5590
function          6455
industry           4983
employment_type     3771
company_profile     3308
requirements        2696
location            346
description          1
dtype: int64

```

Figure 3: Missing values

Possible Anomalies:

- salary_range has very few entries (only 2,868 out of 17,880), which might indicate companies often don't disclose salary.
- Empty or vague descriptions could be a sign of fraudulent postings.

- department field has a large proportion of missing values, suggesting it might not be commonly used.

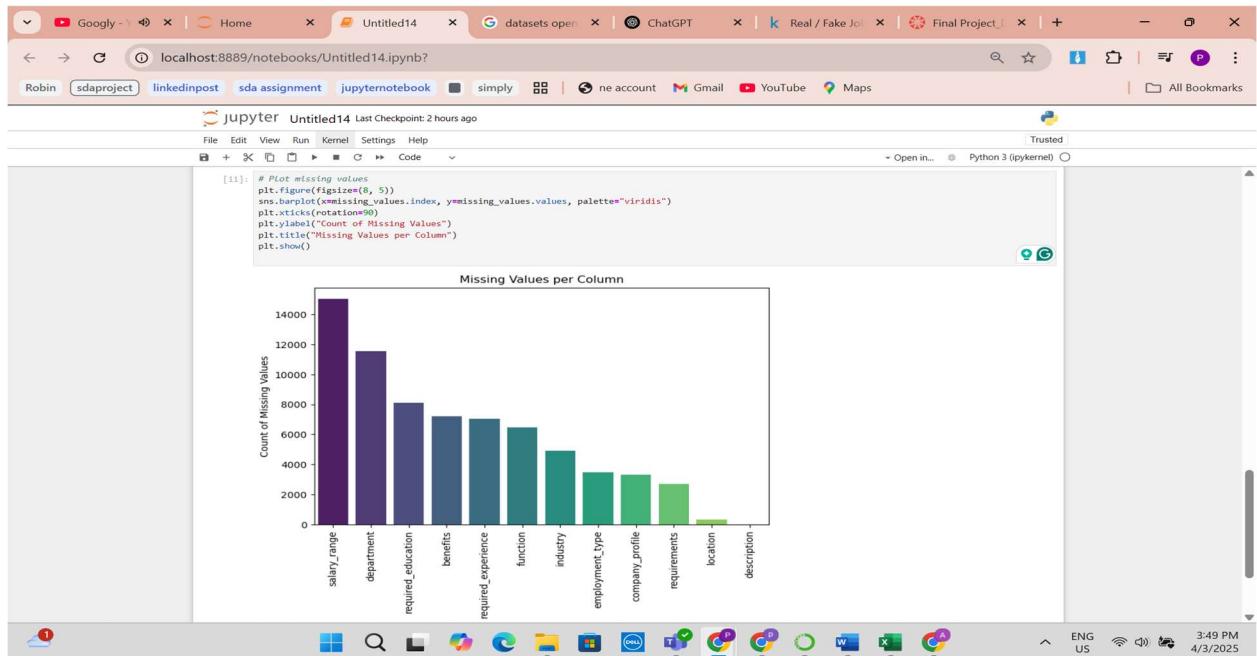


Figure 4: Possible Anomalies

The bar chart highlights the extent of missing values across various columns. The **department**, **salary_range**, and **required_education** fields have significant gaps, which may impact analysis.

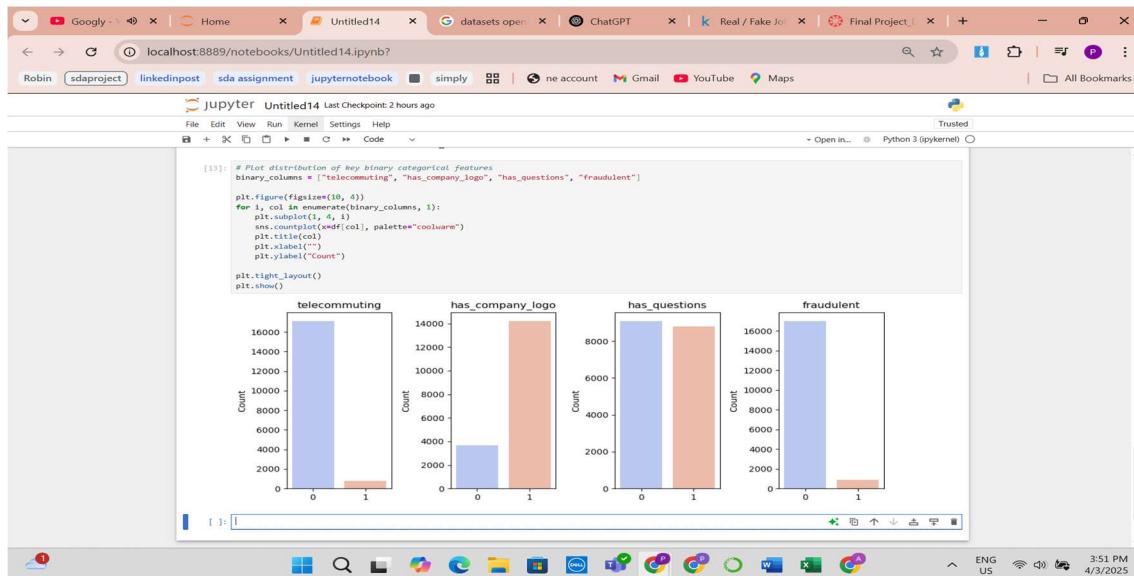


Figure 5: Possible Anomalies

The count plots show some interesting patterns:

- Most job postings do not allow telecommuting (remote work).
- Nearly all job postings have a company logo, suggesting that fraudsters might also include logos to appear legitimate.
- Most postings do not have screening questions, which could indicate either a lack of thorough vetting or a preference for a simpler application process.
- There are significantly fewer fraudulent postings compared to legitimate ones, suggesting an imbalanced dataset.

In the binary categorical feature plots (telecommuting, has_company_logo, has_questions, fraudulent), the blue and orange colors represent the two possible values (0 and 1). Specifically:

- Blue usually represents the count of 0 (False/No).
- Orange represents the count of 1 (True/Yes).

For example, in the fraudulent column:

- Blue represents legitimate job postings (0).
- Orange represents fraudulent job postings (1).

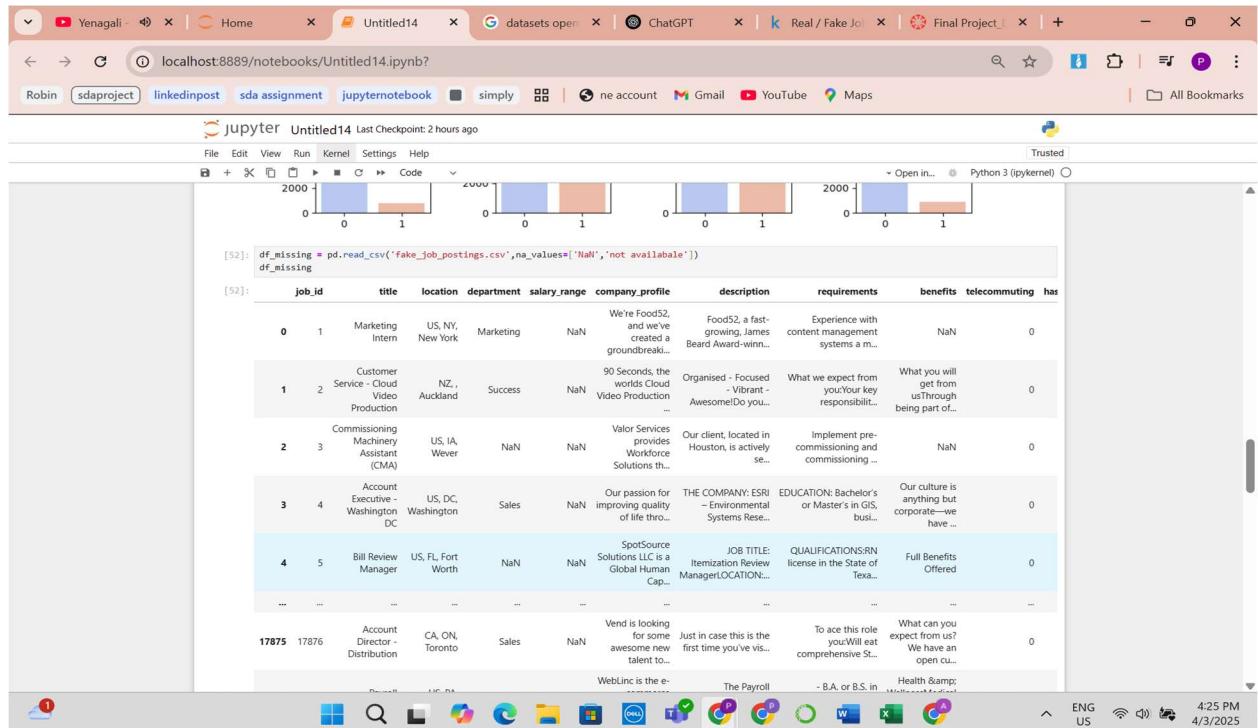


Figure 6: Missing Values

The screenshot shows a Jupyter Notebook interface with the following code:

```

[53]: df_missing.dtypes
[54]: job_id          int64
title           object
location        object
department      object
salary_range    object
company_profile object
description     object
requirements    object
benefits       object
telecommuting   int64
has_logo         int64
has_questions   int64
employment_type object
required_expense object
required_education object
industry        object
forbidden       object
fraudulent      int64
dtype: object

[54]: # Fill missing values with 0
df.fillna(0, inplace=True)

[56]: # Remove rows with missing values
df.dropna(inplace=True)

[57]: df

```

Below the code, the notebook displays a table with 2 rows of data:

	job_id	title	location	department	salary_range	company_profile	description	requirements	benefits	telecommuting	has
0	1	Marketing Intern	US, NY, New York	Marketing	0	We're Food52, and we've created a groundbreak...	Food52, a fast-growing, James Beard Award-winn...	Experience with content management systems a m...	0	0	
1	2	Customer Service - Cloud Video Production	NZ, Auckland	Success	0	90 Seconds, the world's Cloud Video Production	Organised - Focused - Vibrant - Awesome! Do you...	What we expect from you! Your key responsibilit...	What you will get from us through being part of...	0	

Figure 7: Filling Null Values

The screenshot shows a Jupyter Notebook interface with the following code:

```

[56]: # Remove rows with missing values
df.dropna(inplace=True)

[57]: df

```

Below the code, the notebook displays a table with 17876 rows of data:

	job_id	title	location	department	salary_range	company_profile	description	requirements	benefits	telecommuting	has
0	1	Marketing Intern	US, NY, New York	Marketing	0	We're Food52, and we've created a groundbreak...	Food52, a fast-growing, James Beard Award-winn...	Experience with content management systems a m...	0	0	
1	2	Customer Service - Cloud Video Production	NZ, Auckland	Success	0	90 Seconds, the world's Cloud Video Production	Organised - Focused - Vibrant - Awesome! Do you...	What we expect from you! Your key responsibilit...	What you will get from us through being part of...	0	
2	3	Commissioning Machinery Assistant (CMA)	US, IA, Wever	0	0	Valor Services provides workforce solutions th...	Our client, located in Houston, is actively se...	Implement pre-commissioning and commissioning ...	0	0	
3	4	Account Executive - Washington DC	US, DC, Washington	Sales	0	Our passion for improving quality of life thro...	THE COMPANY: ESRI - Environmental Systems Rese...	EDUCATION: Bachelor's or Master's in GIS, busi...	Our culture is anything but corporate—we have ...	0	
4	5	Bill Review Manager	US, FL, Fort Worth	0	0	SpiceSource Solutions LLC is a Global Human Cap...	JOB TITLE: Iteration Review Manager, LOCATION:...	QUALIFICATIONS: RN license in the State of Texa...	Full Benefits Offered	0	
...	
17875	17876	Account Director - Distribution	CA, ON, Toronto	Sales	0	Vend is looking for some awesome new talent to...	Just in case this is the first time you've visi...	To ace this role you'll eat comprehensive St...	What can I expect from you? Will eat healthy cu...	0	
17876	17877	Payroll Accountant	US, PA, Philadelphia	Accounting	0	WebLine is the e-commerce platform and	The Payroll Accountant will focus primarily on...	- B.A. or B.S. in Accounting - Desire to have t...	Health & Wellness/Medical Plan/Prescription	0	

Figure 8: Dropping Values

A screenshot of a Jupyter Notebook interface. The notebook has a single cell containing the following Python code:

```
[68]: print(df.isnull().sum())

```

The output of the code is displayed below the cell, showing the count of null values for each column in the DataFrame:

Column	Count
job_id	0
title	0
location	0
department	0
salary_range	0
company_profile	0
description	0
requirements	0
benefits	0
telecommuting	0
has_company_logo	0
has_questions	0
employment_type	0
required_experience	0
required_education	0
industry	0
function	0
fraudulent	0
dtype: int64	

Figure 9: Checking Null Values

Outlier and Anomalies detection and removal:

Why It Matters

Anomaly detection is used in:

- Fraud detection
- Cybersecurity (intrusion detection)
- Healthcare (rare diseases)
- Manufacturing (defect detection)

A screenshot of a Jupyter Notebook interface. The notebook has two cells. The first cell contains Python code for outlier detection and removal using the IQR method:

```
[70]: # Outlier Detection and Removal using IQR
numerical_cols = ["telecommuting", "has_company_logo", "has_questions", "fraudulent"]
for col in numerical_cols:
    Q1 = df[col].quantile(0.25)
    Q3 = df[col].quantile(0.75)
    IQR = Q3 - Q1
    lower_bound = Q1 - 1.5 * IQR
    upper_bound = Q3 + 1.5 * IQR
    df = df[(df[col] >= lower_bound) & (df[col] <= upper_bound)]
```

The second cell displays the resulting DataFrame, which has been cleaned of outliers. The DataFrame has columns: job_id, title, location, department, salary_range, company_profile, description, requirements, benefits, telecommuting, and n.

Index	job_id	title	location	department	salary_range	company_profile	description	requirements	benefits	telecommuting	n
0	1	Marketing Intern	US, NY, New York	Marketing	0	We're Food52, a fast-growing, James Beard Award-winn...	Experience with content management systems a must...	0	0		
1	2	Customer Service - Cloud Video Production	NZ, Auckland	Success	0	90 Seconds, the world's first video production	Organized - Focused - Vibrant - Awesome! Do you...	What we expect from you! Your key responsibility...	0		
2	3	Commissioning Machinery Assistant (CMA)	US, IA, Wever	0	0	Value services provides Workforce Solutions th...	Our client, located in Houston, is actively see...	Implement pre-commissioning and commissioning ...	0	0	
3	4	Account Executive - Washington DC	US, DC, Washington	Sales	0	Our passion for improving quality of life of th...	THE COMPANY: ESR - Environmental Systems Rese...	EDUCATION: Bachelor's or Master's in GIS, busin...	0		
4	5	Bill Review Manager	US, FL, Fort Worth	0	0	SpotSource Solutions LLC is a Global Human Cap...	JOB TITLE: Itemization Review Manager LOCATION:...	QUALIFICATIONS: RN license in the State of Texa...	Full Benefits Offered	0	

Figure 10: Data Cleaning

Jupyter Untitled14 Last Checkpoint: 2 hours ago

File Edit View Run Kernel Settings Help Trusted

[72]: df

	job_id	title	location	department	salary_range	company_profile	description	requirements	benefits	telecommuting	h
0	1	Marketing Intern	US, NY, New York	Marketing	0	We're Food52, and we've created a groundbreak...	Food52, a fast-growing, James Beard Award-winn...	Experience with content management systems a m...	0	0	
1	2	Customer Service - Cloud Video Production	NZ., Auckland	Success	0	90 Seconds, the world's Cloud Video Production ...	Organized - Focused - Vibrant - Awesome! Do you...	What we expect from you? Your key responsibilit...	What you will get from us through being part of...	0	
2	3	Commissioning Machinery Assistant (CMA)	US, IA, Wever	0	0	Valor Services provides Workforce Solutions th...	Our client, located in Houston, is actively see...	Implement pre-commissioning and commissioning ...	0	0	
3	4	Account Executive - Washington DC	US, DC, Washington	Sales	0	Our passion for improving quality of life thro...	THE COMPANY: ESR - Environmental Systems Rese...	EDUCATION: Bachelor's or Master's in GIS, busi...	Our culture is anything but corporate—we have ...	0	
4	5	Bill Review Manager	US, FL, Fort Worth	0	0	SpotSource Solutions LLC is a Global Human Cap...	JOB TITLE: Iteration Review Manager LOCATION: ...	QUALIFICATIONS: RN license in the State of Tex...	Full Benefits Offered	0	
...	
17872	17873	Product Manager	US, CA, San Francisco	Product Development	0	Flite delivers ad innovation at scale to the w...	Flite's SaaS display ad platform fuels the wor...	BA/BS in Computer Science or a related technic...	Competitive base + attractive stock option plan / Me...	0	
17873	17874	Recruiting Coordinator	US, NC, Charlotte	0	0	0	RESPONSIBILITIES: Will facilitate the recruit...	REQUIRED SKILLS: Associates Degree or a combin...	0	0	
17875	17876	Account Director - Distribution	CA, ON, Toronto	Sales	0	Vend is looking for some awesome new	Just in case this is the first time you've vis...	To ace this role you will eat comprehensive St...	What can you expect from us? We have an	0	

Figure 11: Outlier Detection and Removal

Outlier plots:

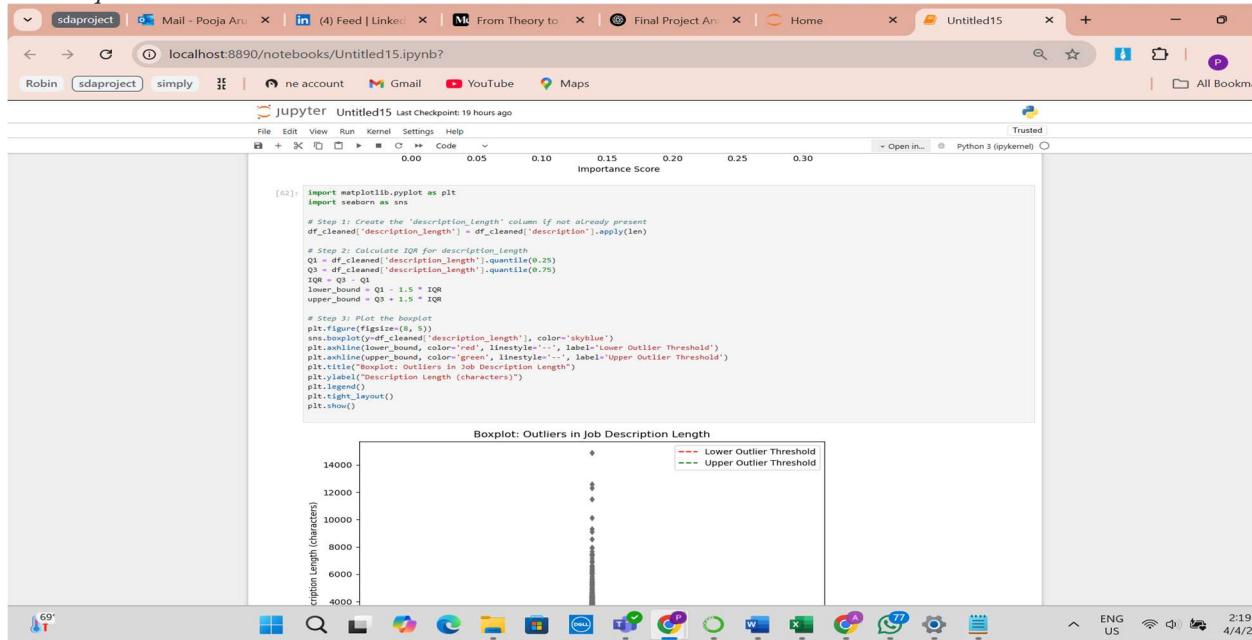


Figure 12: Outlier Detection and Removal

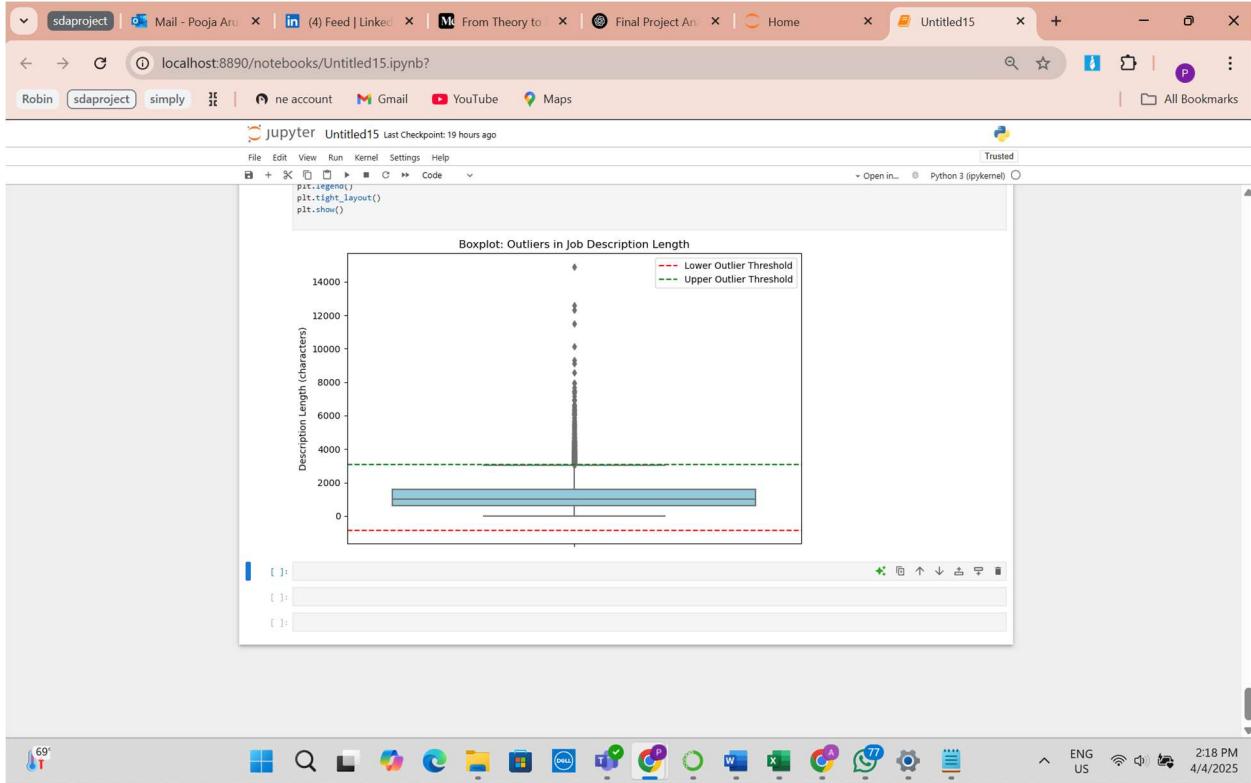


Figure 13:Outlier Detection and Removal

Boxplot Insight – Description Length Outliers

- The box represents the middle 50% of job descriptions:
 - These typically fall between ~500 and ~1500 characters, meaning most job posts have reasonable-length descriptions.
- The green line marks the upper threshold for normal values.
- The red line marks the lower threshold.
- The black dots above the green line are outliers — job descriptions that are unusually long, in some cases reaching over 14,000 characters!

Why this matters:

- Extremely short descriptions (below lower bound) could indicate:
 - Fake/scam listings with little effort put into detail
 - Attempt to avoid detection or copy/paste scams
- Extremely long descriptions (above upper bound) might mean:
 - Overcompensating for legitimacy
 - Spam or repeated irrelevant content
 - Copied content from multiple listings

Conclusion:

Outliers in description_length are valuable indicators of **irregular or suspicious behavior** in job postings. This supports the use of this feature in your fraud detection model and shows why it's important to visualize and monitor such anomalies.

2) Violin Plot Analysis: Job Description Length vs Fraudulent Label

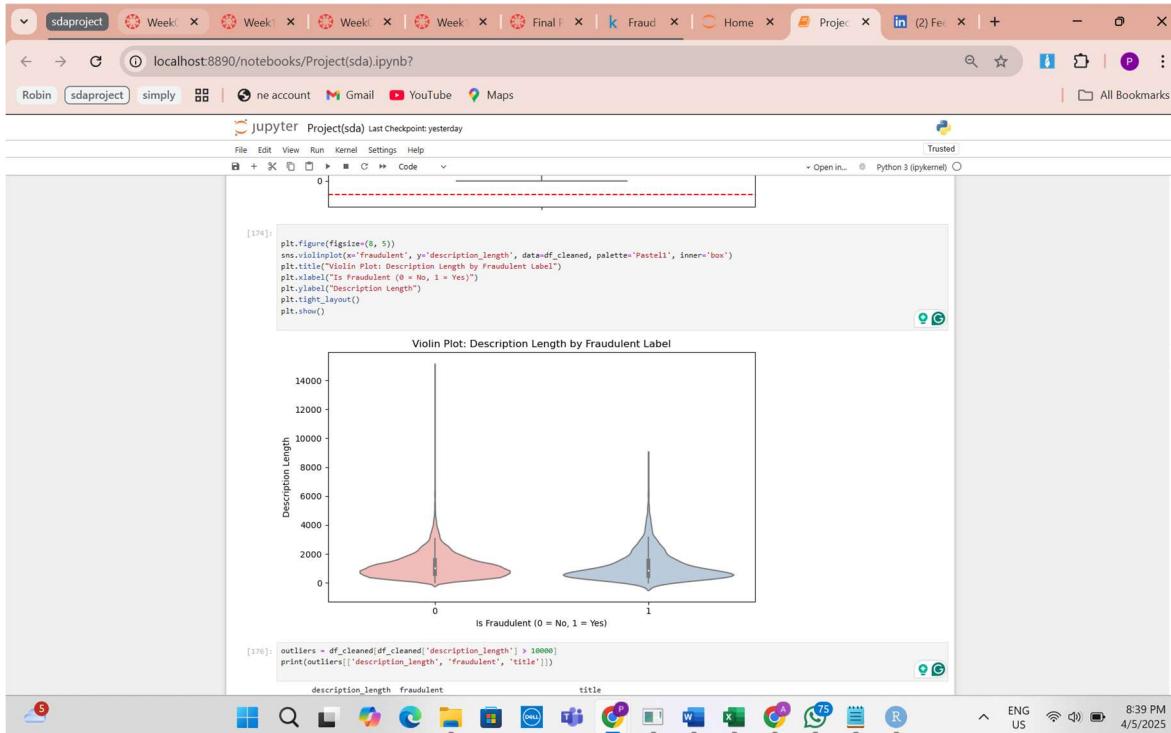


Figure 14: Outlier Detection and Removal

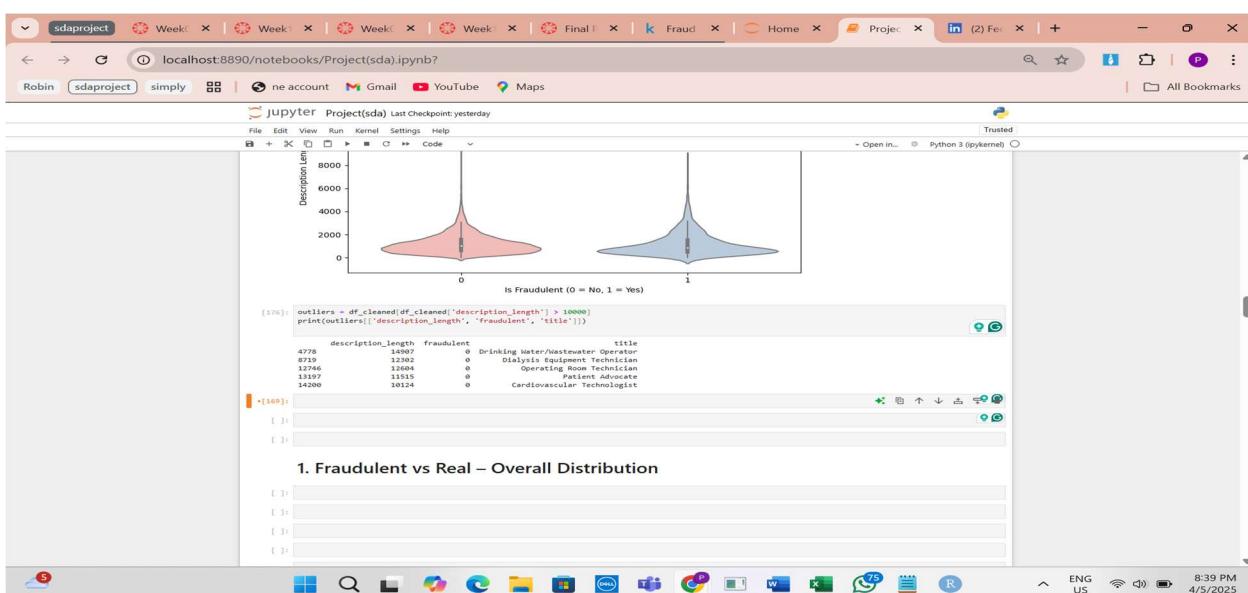


Figure 15: Outlier Detection and Removal

To visually inspect potential outliers in job descriptions, a violin plot was generated comparing the `description_length` distribution between fraudulent and non-fraudulent job postings.

The plot combines a boxplot and a kernel density estimate, allowing us to observe not only the central tendency and spread of each group but also the shape of their distribution.

Key Observations:

- Real job postings (`fraudulent = 0`) show a wide distribution, with most description lengths clustering between 500 to 2000 characters. However, there is a notable long tail extending beyond 14,000 characters, suggesting the presence of extreme outliers — potentially overly detailed or artificially inflated descriptions.
- Fraudulent job postings (`fraudulent = 1`) have a tighter distribution overall, mostly between 200 to 1500 characters, but still exhibit a few unusually long descriptions, indicating that some fake listings may attempt to appear legitimate through verbose or copied content.
- In both cases, the narrow vertical extensions (violin "necks") at the higher end of the y-axis highlight rare and extreme values — a clear visual signal of outliers.

Conclusion:

The violin plot effectively reveals outliers in job description length. These outliers may be important indicators in identifying patterns of deception in job postings. For example, extremely long descriptions in fake jobs may be a tactic to distract or mislead applicants. Conversely, legitimate jobs with unusually long postings may warrant manual review to validate quality and consistency.

Effect of Outlier / anomalies Removal:

To understand the impact of outliers in the `description_length` feature, we generated violin plots comparing the distribution of real and fake job descriptions before and after IQR-based outlier removal.

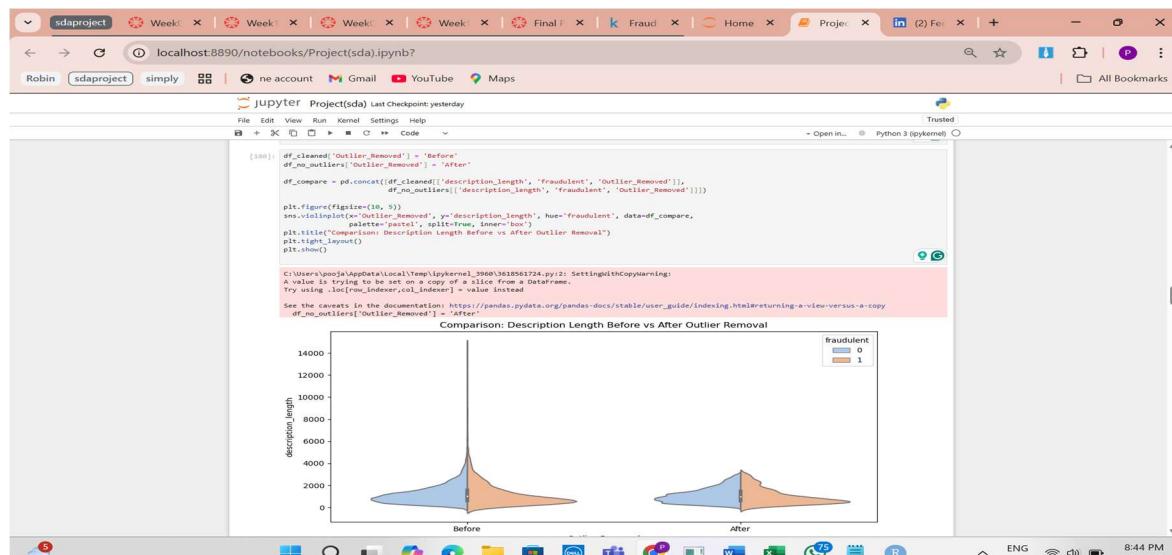


Figure 16: Outlier Detection and Removal effect

What the Visualization Shows:

- Before outlier removal, the distribution for real jobs (`fraudulent = 0`) had a long upper tail stretching past 14,000 characters. This indicates the presence of extreme outliers — job postings with unusually long descriptions that skewed the overall distribution.
- For fake jobs (`fraudulent = 1`), the spread was already narrower but still included a few high-end outliers, suggesting some fraudulent postings may be artificially verbose.
- After removing outliers using the IQR method:
 - The distribution became tighter and more symmetrical for both classes.
 - The extreme long tails disappeared, making it easier to compare the core patterns of description lengths between real and fake job postings.
 - This resulted in a more robust and interpretable view of the typical behavior in the dataset.

Conclusion:

Outlier removal significantly reduced skewness in `description_length`, allowing for more meaningful comparisons between real and fake job descriptions. This step is especially useful before modeling, as it prevents extreme values from dominating the learning process. However, since these outliers may carry signals of fraud, they should also be examined separately for deeper insight.

2)

The screenshot shows a Jupyter Notebook interface with multiple tabs at the top, including 'sdaproject', 'Week1', 'Week2', 'Week3', 'Final Project', 'Fraud', 'Home', 'Project', 'LinkedIn', and '(2) Feb'. The main area displays a code cell and its output. The code cell contains:`[189]: function_counts = df_cleaned['function'].value_counts()
print(function_counts)`
The output shows a series of function names and their counts, such as Information Technology (474), Sales (424), Engineering (296), Customer Service (243), Marketing (183), Administrative (155), Other (78), Management (72), Accounting/Auditing (69), Design (64), Project Management (54), Business Development (52), Finance (44), Human Resources (42), Advertising (36), Writing/editing (27), Health Care Provider (27), Quality Assurance (26), Education (25), Data Analyst (24), Art/Create (22), Manufacturing (19), Production (19), General Business (19), Consulting (18), Product Management (18), Business Analyst (14), Training (13), Strategy/Planning (11), Financial Analyst (11), Public Relations (10), Supply Chain (8), Research (5), Dissertation (5), Purchasing (5), Legal (4), Science (1).
Name: count, dtype: int64

Below the code cell, another cell is partially visible with the following content:`+ [190] threshold = 40
Find rare functions`

Figure 17: Anomalies Detection

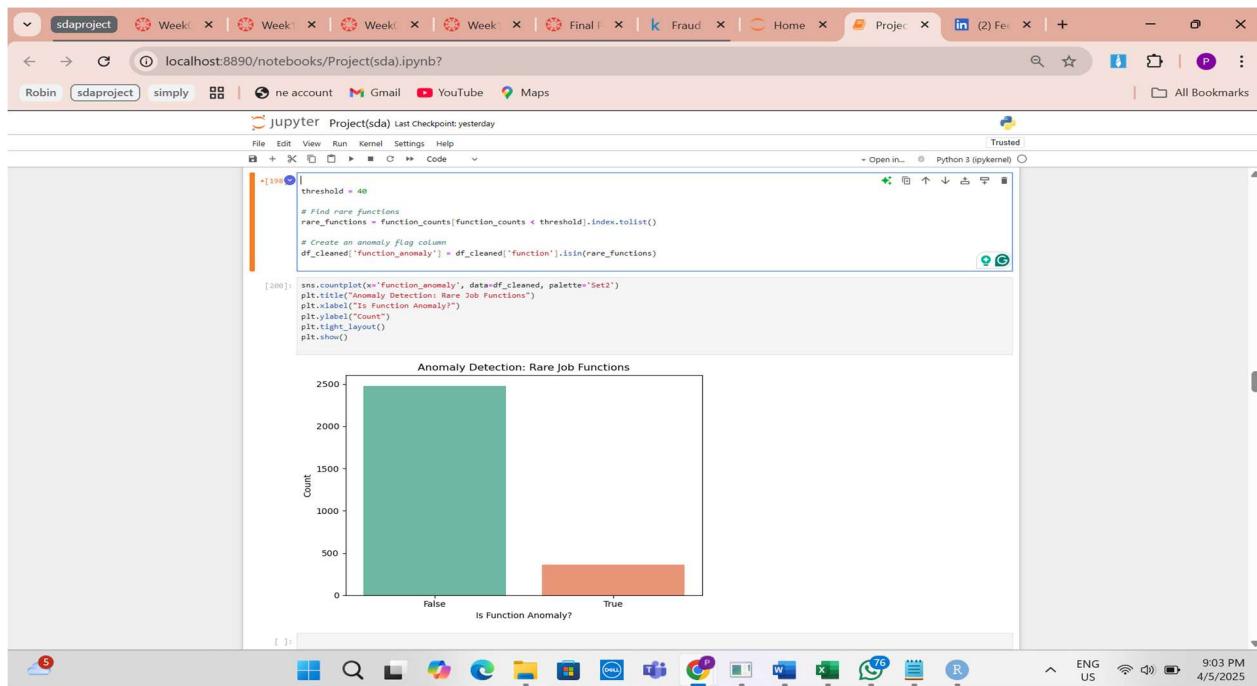


Figure 18: Anomalies Detection

What does this plot show?

Plotted a count of job postings split by:

- `function_anomaly = False` (Normal/common job functions)
- `function_anomaly = True` (Rare/unusual job functions)

With threshold of 40, defined rare functions as those appearing fewer than 40 times in the dataset.

Anomaly Detection in Job Function Column

To identify unusual patterns in job functions, we applied categorical anomaly detection based on frequency. Job functions that appeared less than 40 times in the dataset were flagged as anomalies.

Findings:

- Out of all job postings, approximately 350 entries had rare job functions.
- These rare functions could indicate:
 - Niche or specialized roles
 - Data inconsistencies (misspelled or inconsistent naming)
 - Potential fraudulent postings trying to mask under unfamiliar job titles

Conclusion:

Rare job functions represent a non-negligible portion of the dataset and should be carefully reviewed. These anomalies may hold valuable signals in distinguishing fraudulent job posts from legitimate ones.

EDA:

1) Employment Type as a Distinguishing Factor in Fraudulent Job Postings

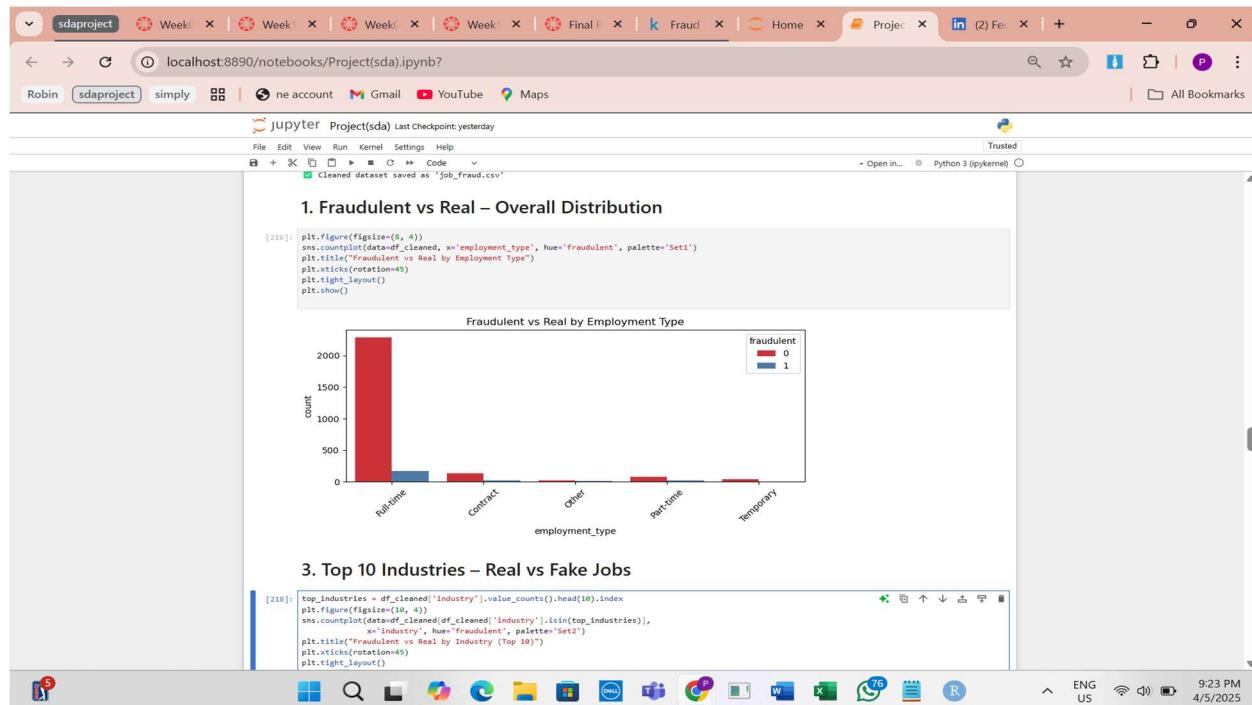


Figure 19: Employment Type

The analysis of employment types across fraudulent and legitimate job postings shows a clear disparity, offering valuable insights into patterns used by fake job listings.

Key Observations:

- Fraudulent job postings are heavily concentrated under the “Full-time” category.
- Legitimate postings are more evenly distributed across various employment types, including Contract, Part-time, and Temporary.
- This concentration suggests that fraudsters intentionally label fake jobs as full-time to increase credibility and appeal to a broader audience.

Why This Matters:

- Full-time positions appear more legitimate, making them an easy tool for fraudsters to gain trust quickly.
- Real jobs reflect greater variety, aligning with actual hiring trends and flexibility in employment structures.

- Monitoring employment type can help identify red flags during the early stages of detection.

Implications for Detection:

- The employment_type field can serve as a predictive feature in identifying fraudulent job posts.
- When combined with other signals (e.g., vague descriptions, missing logos, unrealistic salaries), full-time labels on questionable postings may increase the likelihood of fraud.
- This reinforces the value of using data analysis techniques to uncover behavioral patterns that are not immediately obvious through manual review.

2) Industry Distribution in Real vs Fraudulent Job Postings

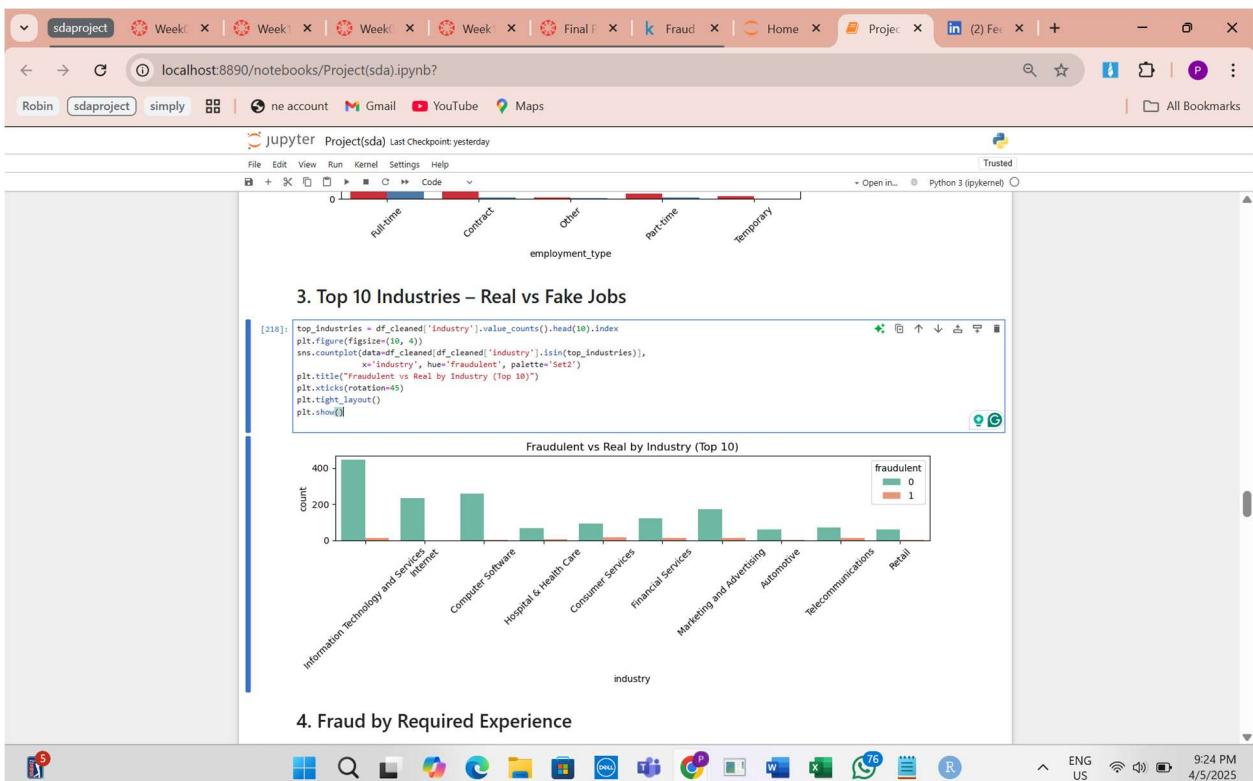


Figure 20: Top 10 Industries

Visualization compares the distribution of fake and real job postings across the top 10 most common industries in the dataset. The findings reveal distinct differences in how fraudsters target specific sectors.

Key Observations:

- Industries like Information Technology & Services, Internet, and Computer Software have the highest number of legitimate job postings — and also the most fraud activity among all categories.

- Fraudulent postings are relatively less common in industries such as Retail, Telecommunications, and Automotive.
- Highly digital sectors are more exposed to fake jobs, likely because they attract more online applicants and are easier for scammers to mimic.

Why This Matters:

- Fraudsters may target fast-growing, tech-oriented industries to exploit high demand and applicant volume.
- Real jobs show broad diversity across industries, whereas fake postings are concentrated in specific high-tech sectors, making industry type a useful predictive feature.
- Identifying industry-based risk zones can help in prioritizing fraud detection efforts or implementing stricter screening in high-risk sectors.

Implications for Detection:

- Industry is a valuable feature in classification models used for detecting fake job postings.
- Jobs in industries like IT, Internet, and Software should be examined more closely when paired with other risk factors.
- Platforms and job boards can leverage this information to allocate risk scoring to specific industries for enhanced security.

3) Required Experience in Real vs Fraudulent Job Postings:

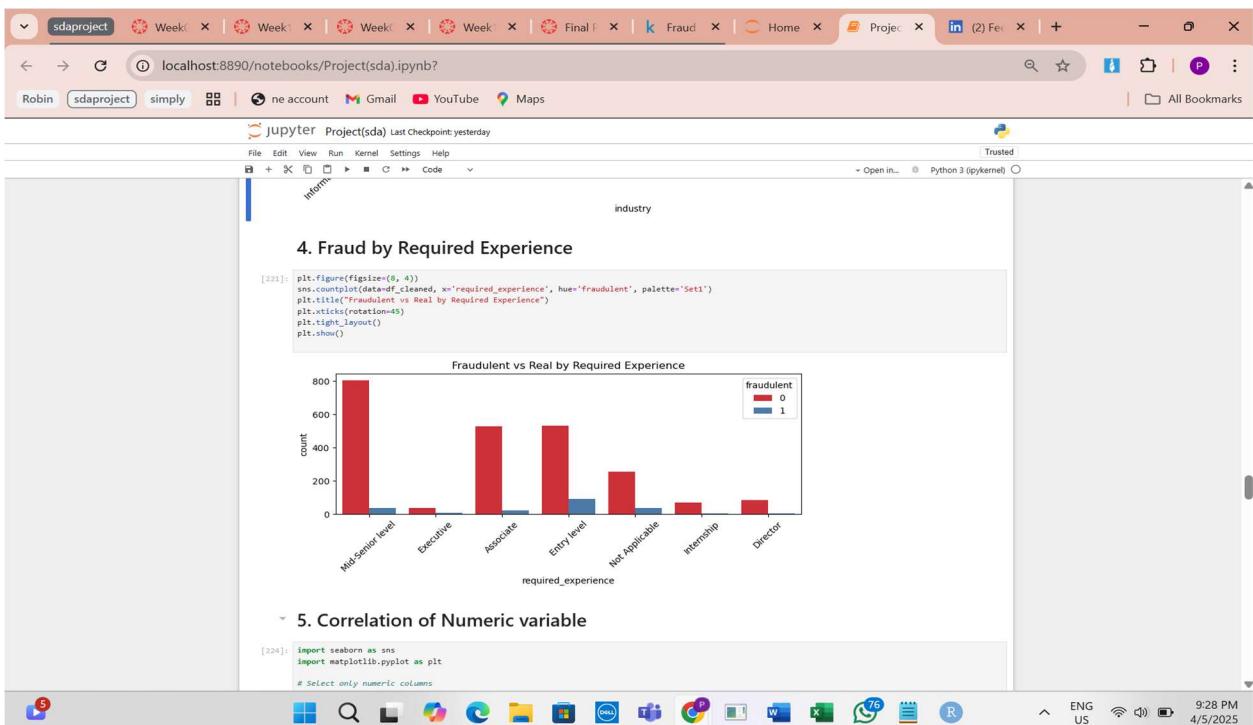


Figure 21: Fraud By required experience field

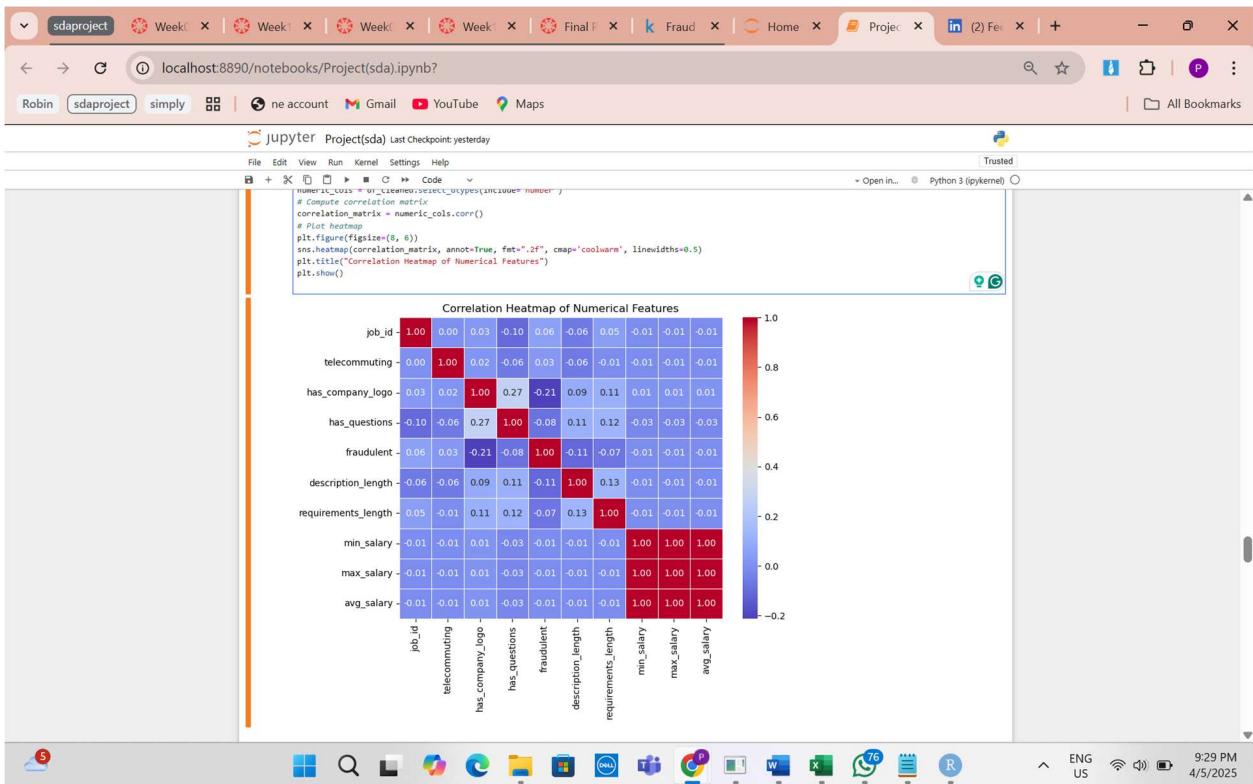


Figure 22: Correlation Heatmap

This heatmap visualizes the relationships between numerical variables in the dataset, particularly focusing on how each correlates with the target variable `fraudulent`.

Key Observations:

- The variable `has_company_logo` shows a negative correlation (-0.21) with `fraudulent`, suggesting that fraudulent job postings are more likely to lack a company logo.
- `has_questions` has a positive correlation (0.27) with `fraudulent`, implying that fraudulent jobs often include screening questions, possibly to seem more legitimate.
- Other features like `description_length`, `requirements_length`, and salary-related fields show very weak correlations with `fraudulent` — indicating that length or salary alone may not be strong individual predictors.

Why This Matters:

- Identifying correlations helps prioritize which features to focus on when building predictive models.
- Features like presence of a company logo and inclusion of screening questions could be strong flags in detecting fraudulent listings.
- Low correlation doesn't always mean useless — it may still matter in multivariate models (like random forests), but is less informative in isolation.

Implications for Detection:

- Correlation analysis confirms that certain binary flags (like logo presence or questions) are statistically linked to fraud and can be used to enhance model performance.
- While some features (like salary and description length) are important for analysis, their direct correlation with fraud is low — suggesting they are supporting signals rather than primary indicators.
- This analysis ensures that model-building focuses on the most informative features for detecting anomalies and fraud patterns.

5) Employment Type Distribution Across Top Industries

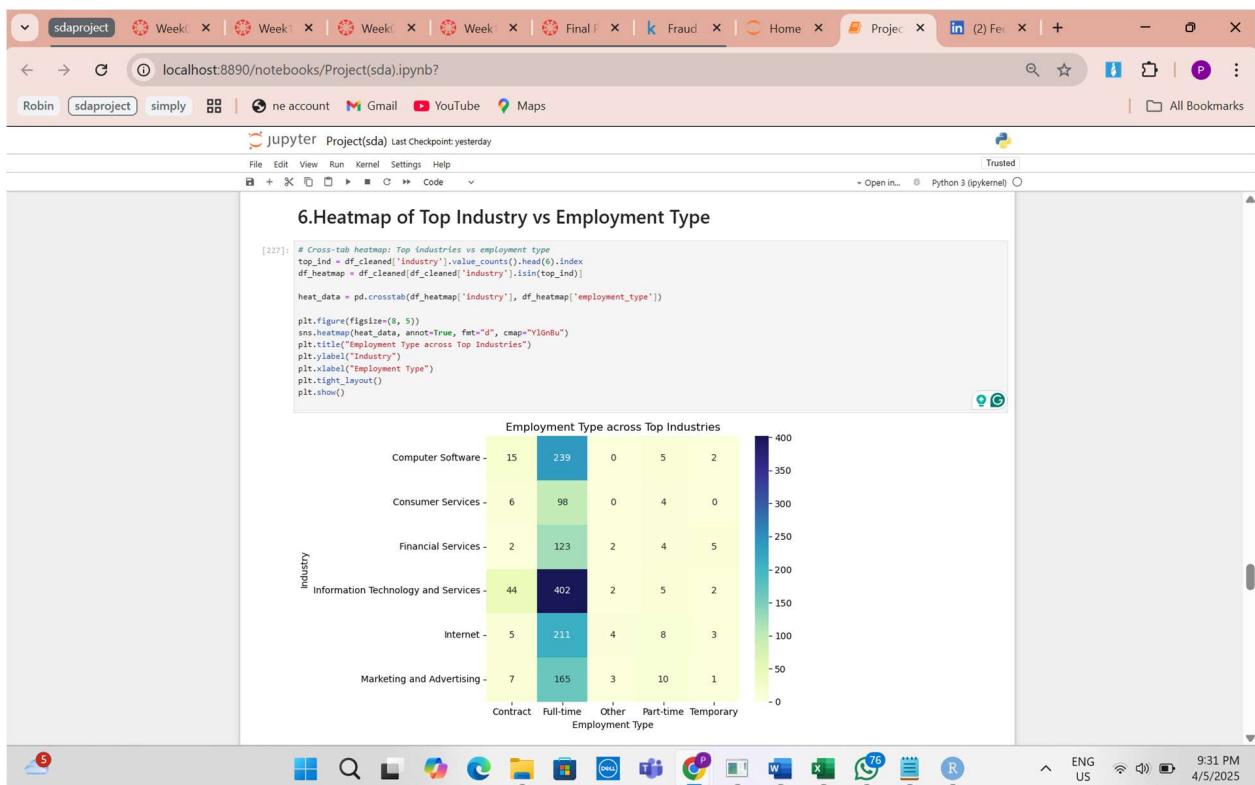


Figure 23: Employment time across Industries

This heatmap shows the cross-distribution of employment types within the top six most common industries in the dataset. It highlights how different industries prefer or rely on specific types of employment contracts.

Key Observations:

- Full-time employment is the most dominant type across all top industries.
- Information Technology and Services and Computer Software have the highest concentration of full-time roles, with over 400 and 200 postings respectively.
- Contract roles are minimal in comparison, especially in industries like Internet and Consumer Services.

- Employment diversity (e.g., inclusion of part-time, temporary) is limited in most industries, suggesting a hiring preference toward long-term full-time positions.

Why This Matters:

- Industries with a high share of full-time roles may also be more vulnerable to fraudulent postings, as scammers tend to mimic standard hiring practices to appear trustworthy.
- Understanding how employment types vary by industry helps build context-aware fraud detection systems — for instance, a part-time job in an industry that almost always hires full-time might be worth flagging.

Implications for Analysis & Detection:

- When combined with earlier fraud detection findings, this heatmap strengthens the case for using employment type + industry as a joint feature in predictive models.
- Industries that post mostly full-time jobs could require additional validation layers to ensure job authenticity.
- This cross-tab analysis improves targeted monitoring, allowing fraud detection efforts to focus on the most vulnerable industry-employment combinations.

6) Common Words in Real Job Descriptions

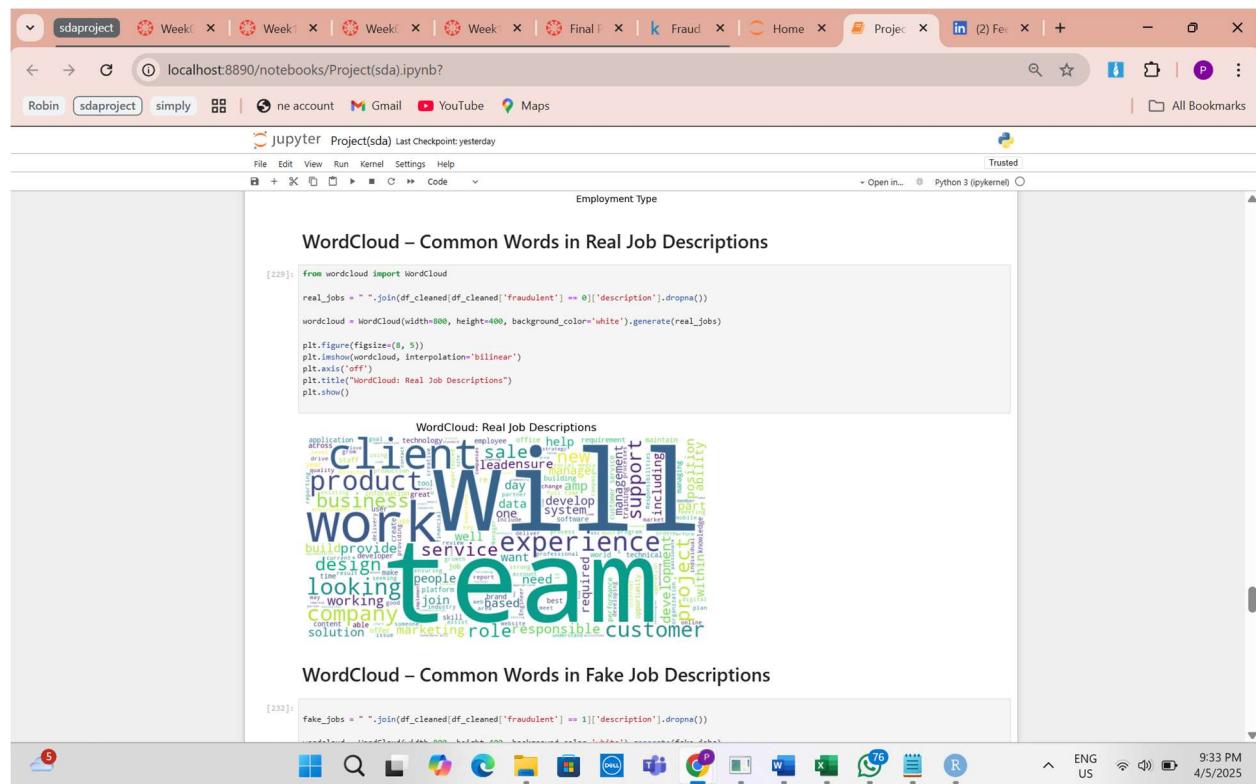


Figure 24: Word Cloud (Common Words in Real Jobs)

This WordCloud highlights the most frequently used words in real (non-fraudulent) job postings, revealing key characteristics of legitimate listings.

Key Observations:

- Prominent words include “team,” “client,” “work,” “experience,” “service,” and “company.”
- The language strongly emphasizes collaboration, service delivery, and professional experience.
- Terms like “looking,” “role,” and “responsible” suggest a formal hiring structure where expectations are clearly defined.

Why This Matters:

- Real job descriptions focus on organizational fit, job responsibilities, and required qualifications.
- They use structured, professional language that reflects genuine business needs and expectations.
- These words align with standard corporate terminology and suggest clear job objectives and deliverables.

Implications for Detection:

- The presence of such words can be used as positive indicators in distinguishing real job descriptions.
- Language models or text classifiers can be trained to recognize patterns aligned with legitimate job postings, increasing accuracy in fraud detection.
- This baseline helps set up a clear contrast when analyzing the fake job descriptions next.

7) Common Words in Fake Job Descriptions:

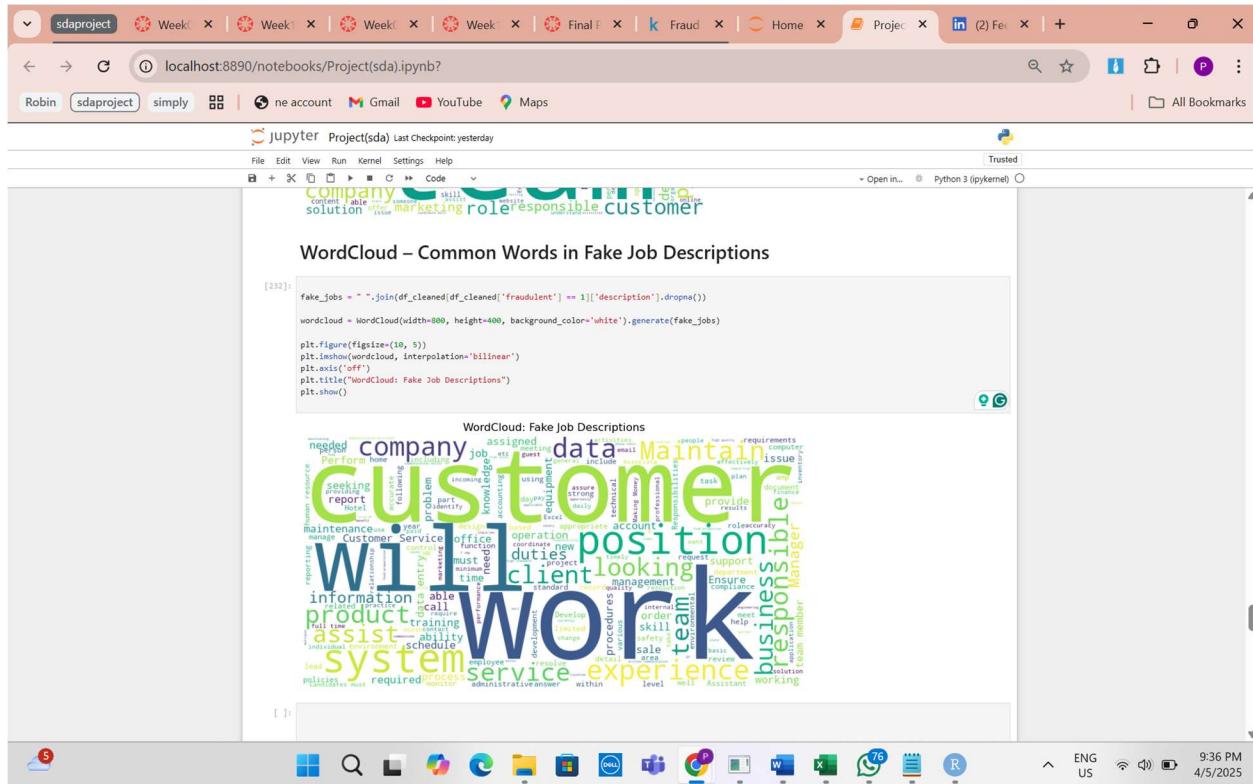


Figure 25:Common Words in Fake Jobs

The WordCloud visualization for fraudulent job postings displays the most frequently used words in fake job descriptions. These words help uncover the linguistic patterns scammers rely on to attract candidates.

Key Observations:

- Prominent words include “customer,” “work,” “system,” “position,” “maintain,” “assist,” and “product.”
- Many of these words are generic, vague, or task-oriented, lacking the contextual richness seen in real job descriptions.
- The repeated use of terms like “assist,” “maintain,” “call,” and “required” reflects basic role language, often without specific responsibilities or qualifications.

Why This Matters:

- Fake job descriptions tend to use broad, ambiguous terms to appeal to a wide audience and avoid scrutiny.
- Unlike real jobs, which focus on teamwork, goals, and qualifications, fake postings lack personal, company-specific, or outcome-driven language.
- This generic phrasing is likely intentional, allowing scammers to copy-paste across different roles and platforms.

Comparison with Real Job Descriptions:

Real Jobs Emphasize	Fake Jobs Emphasize
Team collaboration	Basic tasks
Experience and client handling	System maintenance, customer
Specific goals or platforms	Vague roles like "assist"
Role clarity and responsibility	Generic and repeated phrases

Implications for Detection:

- The **absence of context-rich keywords** can be a strong indicator of fraud.
- NLP techniques like TF-IDF, keyword frequency, or embedding-based classifiers can be used to **differentiate authentic descriptions from fakes**.
- Text-based patterns are critical for detection, especially when metadata (like salary or logo) is missing.

8) Overall Distribution of Real vs Fake Job Postings

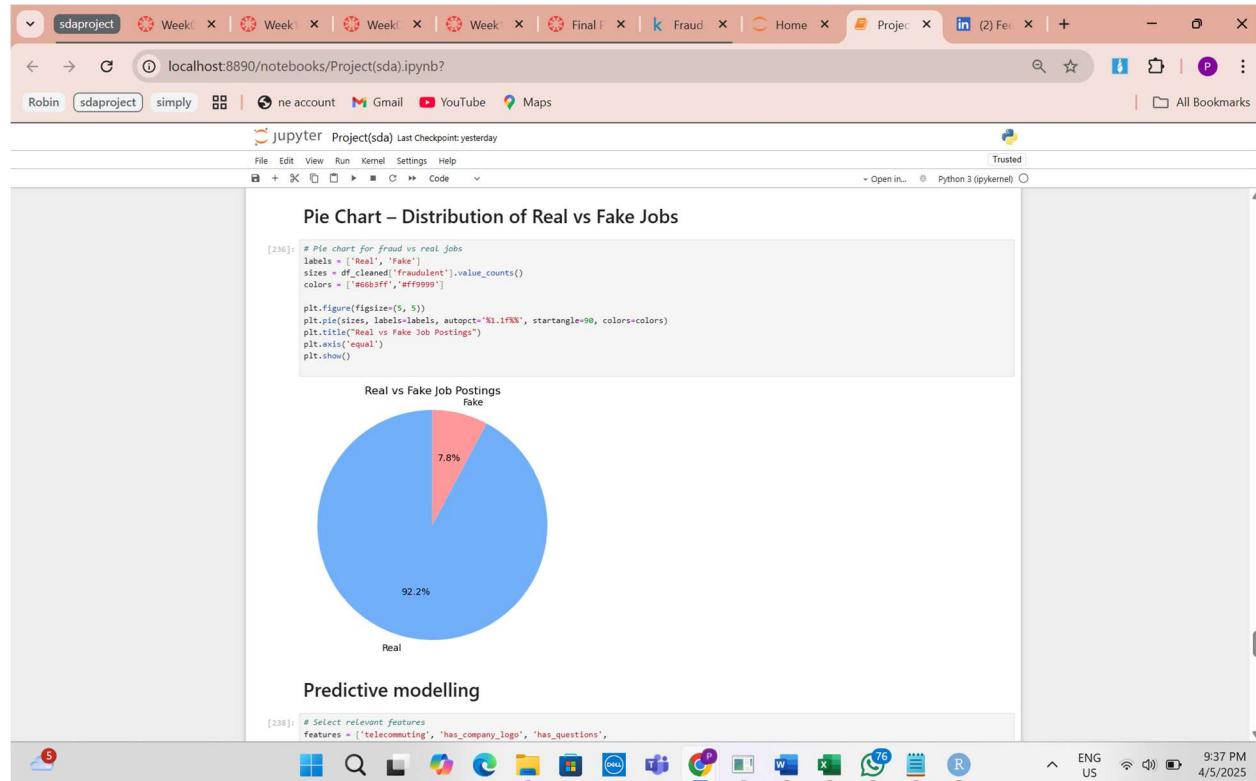


Figure 26: Real vs Fake job distribution

This pie chart presents the class distribution in the dataset, separating real and fake job listings based on the fraudulent label.

Key Observations:

- Real job postings make up 92.2% of the dataset.
- Fake job postings account for just 7.8%, highlighting a significant class imbalance.

Why This Matters:

- The dataset is heavily imbalanced, which is common in fraud detection problems.
- While fake postings represent a small portion, they are critical to detect, as they can cause real-world harm (e.g., scams, identity theft).
- Class imbalance poses a challenge in modeling — traditional accuracy metrics become misleading (e.g., predicting “all real” would still yield >90% accuracy).

Implications for Modeling:

- Balanced accuracy, precision, recall, and F1-score are better metrics for evaluating fraud detection models.
- Techniques such as oversampling (e.g., SMOTE) or undersampling may be necessary to ensure the model pays attention to the minority class (fraud).
- It emphasizes the importance of feature selection, anomaly detection, and careful model evaluation to handle imbalanced data effectively

Predictive Modelling:

Predictive Modeling: Detecting Fraudulent Job Postings

The screenshot shows a Jupyter Notebook interface with the title "Predictive modelling". The code cell contains the following Python script:

```
[238]: # Select relevant features
features = ['telecommuting', 'has_company_logo', 'has_questions',
            'employment_type', 'required_experience', 'required_education',
            'industry', 'function']

target = 'fraudulent'

[239]: from sklearn.preprocessing import LabelEncoder

# Copy dataset and drop rows with 'Unknown' for cleaner training
df_ml = df_cleaned[df_cleaned[features].ne('Unknown').all(axis=1)].copy()

# Encode categorical variables
le = LabelEncoder()
for col in features:
    df_ml[col] = le.fit_transform(df_ml[col])

# Define X and y
X = df_ml[features]
y = df_ml[target]

[240]: from sklearn.model_selection import train_test_split

# 70% train, 30% test
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

[241]: from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, confusion_matrix

# Create model with class_weight to handle imbalance
rf_model = RandomForestClassifier(random_state=42, class_weight='balanced')
rf_model.fit(X_train, y_train)

# Predictions
rf_preds = rf_model.predict(X_test)

# Evaluation
```

Figure 27: Predictive Modelling

The screenshot shows a Jupyter Notebook interface with the title "Random Forest Classifier Report". The code cell contains the following Python script:

```
[242]: # Evaluation
print("Random Forest Classification Report:\n")
print(classification_report(y_test, rf_preds))

print("Confusion Matrix:\n")
print(confusion_matrix(y_test, rf_preds))

Random Forest Classification Report:
      precision    recall  f1-score   support
          0       0.97     0.98     0.98    791
          1       0.76     0.60     0.67     62

accuracy                           0.86
macro avg       0.86     0.79     0.82    853
weighted avg    0.95     0.96     0.95    853

Confusion Matrix:
[[779 12]
 [ 25 37]]

[242]: import pandas as pd
import matplotlib.pyplot as plt

# Plot feature importances
importances = rf_model.feature_importances_
feat_df = pd.DataFrame({'Feature': X.columns, 'Importance': importances})
feat_df = feat_df.sort_values(by='Importance', ascending=False)

# Bar plot
plt.figure(figsize=(8, 5))
plt.barh(feat_df['Feature'], feat_df['Importance'], color='skyblue')
plt.title('Feature Importance - Random Forest')
plt.gca().invert_yaxis()
plt.xlabel('Importance Score')
plt.show()
```

A bar chart titled "Feature Importance - Random Forest" is displayed below the code cell, showing the importance score for each feature.

Figure 28: Random Forest Classifier Report

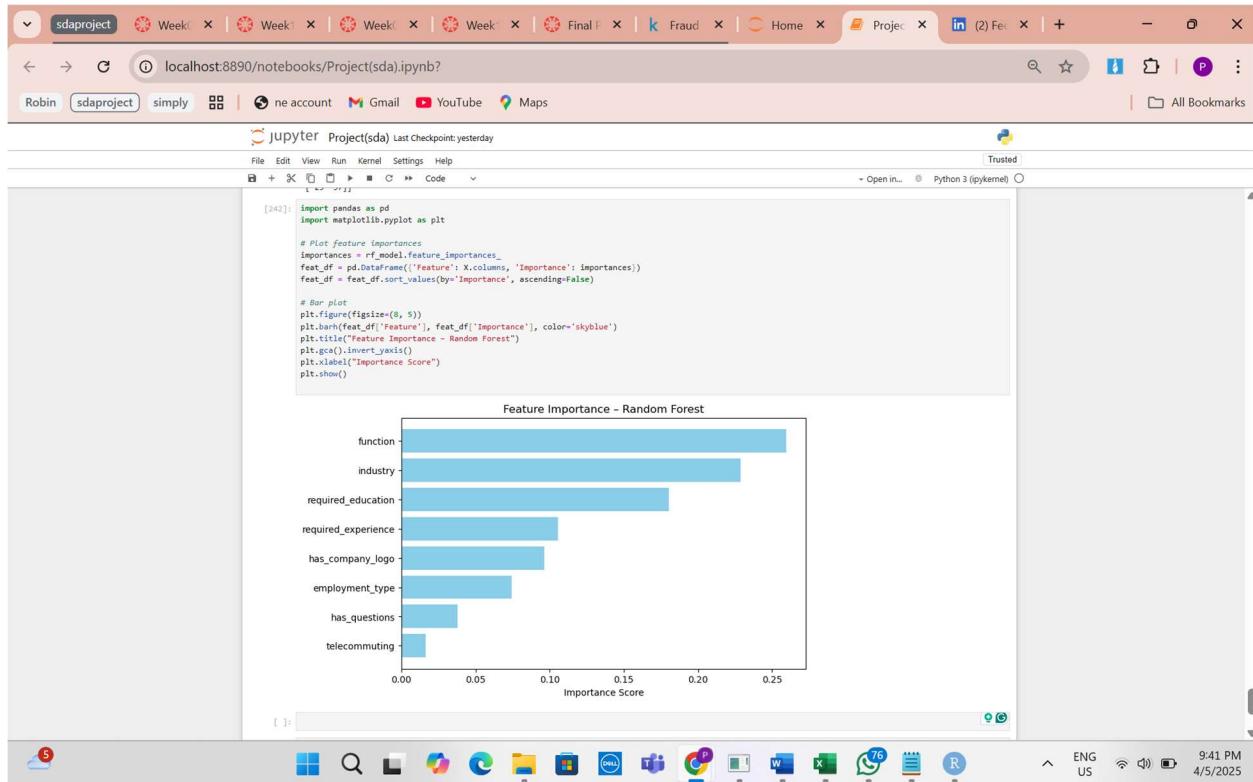


Figure 29: Feature Importance Score

A Random Forest classifier was trained using selected categorical features to predict whether a job posting is fraudulent. The model handled class imbalance using `class_weight='balanced'`, ensuring better sensitivity toward the minority (fraud) class.

Key Model Details:

- Features used: `telecommuting`, `has_company_logo`, `has_questions`, `employment_type`, `required_experience`, `required_education`, `industry`, `function`
- Target: fraudulent (binary: 0 = Real, 1 = Fake)
- Model used: Random Forest Classifier
- Train-test split: 70-30

Model Evaluation Metrics:

Metric	Real Jobs (0)	Fake Jobs (1)
Precision	0.97	0.76
Recall	0.98	0.60
F1-Score	0.98	0.67

- Overall accuracy: 96%
- Macro-average F1-score: 0.82
- Confusion matrix shows 37 fake jobs were correctly identified (true positives), while 25 were missed (false negatives).

Insights:

- The model performs very well on real jobs, and reasonably well on fake jobs, despite class imbalance.
- The use of `class_weight='balanced'` helped improve recall on the minority class, which is essential in fraud detection.
- The F1-score of 0.67 for fraud cases suggests a good balance between catching frauds and avoiding false positives.

Feature Importance Analysis:

The most influential features for detecting fake job postings were:

1. Function – Type of role posted (e.g., Sales, Admin)
2. Industry – The field of the job (e.g., IT, Services)
3. Required Education – Minimum qualifications listed
4. Required Experience – Experience level expected
5. Has Company Logo – Whether a company logo is present

These insights align with your earlier exploratory findings. Scammers often post fake jobs with vague functions, uncommon industries, or minimal details like missing logos.

Pivot table (R studio):

```

1 <-- {r}
2 library(readr) # For reading CSV files
3 df <- read_csv("job_fraud_1.csv")
4 ...
5 ...

Rows: 17879 Columns: 18— column specification

Delimiter: ","
chr (13): title, location, department, salary_range, company_profile, description...
dbl (5): job_id, telecommuting, has_company_logo, has_questions, fraudulent
i Use `spec()` to retrieve the full column specification for this data.
i Specify the column types or set `show_col_types = FALSE` to quiet this message.

6 ...
7 ...
8 library(dplyr)

```

Console Terminal Background Jobs

```

R 4.4.2 - C:/Users/pooja/OneDrive/Documents/structured data analytics/Final Project/FinalProject/
+   company_profile, requirements, location),
+   ~ replace_na(., "Unknown") # Replace NA with "Unknown"
+ )
>
> library(knitr)
>
> df_cleaned %>%
+   group_by(employment_type, fraudulent) %>%
+   summarise(count = n(), .groups = "drop") %>%
+   pivot_wider(names_from = fraudulent, values_from = count, values_fill = 0) %>%
+   rename(Real Jobs = 0, 'Fake Jobs' = 1) %>%
+   kable()
>
> |

```

Figure 30: Read Dataset in R

```

6 ...
7 ...
8 library(dplyr)
9 library(tidyverse)
10
11 df_cleaned <- df %>%
12   distinct() %>%
13   filter(!is.na(description)) %>%
14   mutate(across(
15     c(salary_range, department, required_education, benefits,
16       required_experience, function, industry, employment_type,
17       company_profile, requirements, location),
18     ~ replace_na(., "Unknown")) # Replace NA with "Unknown"
19   )
20 ...
21 ...
22

```

Console Terminal Background Jobs

```

R 4.4.2 - C:/Users/pooja/OneDrive/Documents/structured data analytics/Final Project/FinalProject/
+   company_profile, requirements, location),
+   ~ replace_na(., "Unknown") # Replace NA with "Unknown"
+ )
>
> library(knitr)
>
> df_cleaned %>%
+   group_by(employment_type, fraudulent) %>%
+   summarise(count = n(), .groups = "drop") %>%
+   pivot_wider(names_from = fraudulent, values_from = count, values_fill = 0) %>%
+   rename(Real Jobs = 0, 'Fake Jobs' = 1) %>%
+   kable()
>
> |

```

Figure 31: Data Description

The screenshot shows the RStudio interface with the following components:

- Source View:** Displays R code for creating a pivot table. The code uses dplyr and knitr packages to group by employment type, sum up counts, and then pivot_wider to create columns for 'Real Jobs' and 'Fake Jobs'.
- Console View:** Shows the execution of the R code in the R 4.4.2 environment. The output displays a pivot table titled 'employment_type' with columns 'Real Jobs' and 'Fake Jobs'.
- Environment View:** Shows the global environment with two data frames: 'df' (17879 obs. of 18 variables) and 'df_cleaned' (17879 obs. of 18 variables).
- Files View:** Shows the project structure under 'OneDrive > Documents > structured data analytics > Final Project > FinalProject'. It lists 'FinalProject.Rproj' (267 B, Apr 5, 2025, 3:49 PM) and 'fake_job_postings.csv' (47.7 MB, Apr 3, 2025, 3:16 PM).

Figure 32: Pivot Table for real and Fake Jobs

The screenshot shows the RStudio interface with the following components:

- Source View:** Displays the same R code as in Figure 32.
- Console View:** Shows the execution of the R code. The output displays a pivot table titled 'employment_type' with columns 'Real Jobs' and 'Fake Jobs'.
- Environment View:** Shows the global environment with two data frames: 'df' (17879 obs. of 18 variables) and 'df_cleaned' (17879 obs. of 18 variables).
- Files View:** Shows the project structure under 'OneDrive > Documents > structured data analytics > Final Project > FinalProject'. It lists 'FinalProject.Rproj' (267 B, Apr 5, 2025, 3:49 PM) and 'fake_job_postings.csv' (47.7 MB, Apr 3, 2025, 3:16 PM).

Figure 33: Results for Pivot Table of real and Fake Jobs

Real vs Fake Job Postings by Employment Type

The pivot table reveals important trends in how fraud is distributed across different employment types:

- Full-time jobs dominate the dataset with over 11,000 real listings and 490 fake ones. While this is the most common category, the fraud rate is relatively low (~4%), suggesting these listings are mostly legitimate.
- Part-time jobs show a notably higher fraud rate (~9%), indicating scammers often target people looking for flexible or short-term work.
- "Unknown" employment type also shows a significant fraud risk — with 240 out of 3,470 postings (around 7%) being fake. This suggests that missing or vague employment information may be a red flag.
- Contract and temporary roles have the lowest number of fake listings, indicating these are less frequently targeted by scammers in this dataset.

Conclusion:

The analysis confirms that employment type is a useful feature for identifying fraud. Listings with missing or vague types ("Unknown") and part-time roles show higher fraud rates and should be treated with caution in any detection model.

2)

The screenshot shows the RStudio interface with the following details:

- File Menu:** File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Code Editor:** Untitled1.R contains R code for data cleaning and pivoting. The code includes library imports (dplyr, tidyverse, knitr), data loading, grouping by required experience, summarizing counts, pivoting wider, renaming columns, and arranging results. It ends with a kable() call.
- Data Environment:** Shows two datasets: df (17879 obs. of 18 variables) and df_cleaned (17879 obs. of 18 variables).
- File Explorer:** Shows files in the project directory: FinalProject.Rproj and fake_job_postings.csv.
- Console:** Displays the R session history with the same commands as the code editor.
- System Tray:** Shows battery level, network status, and system time (4:07 PM, 4/5/2025).

Figure 34: Pivot Table Based on Experience Required

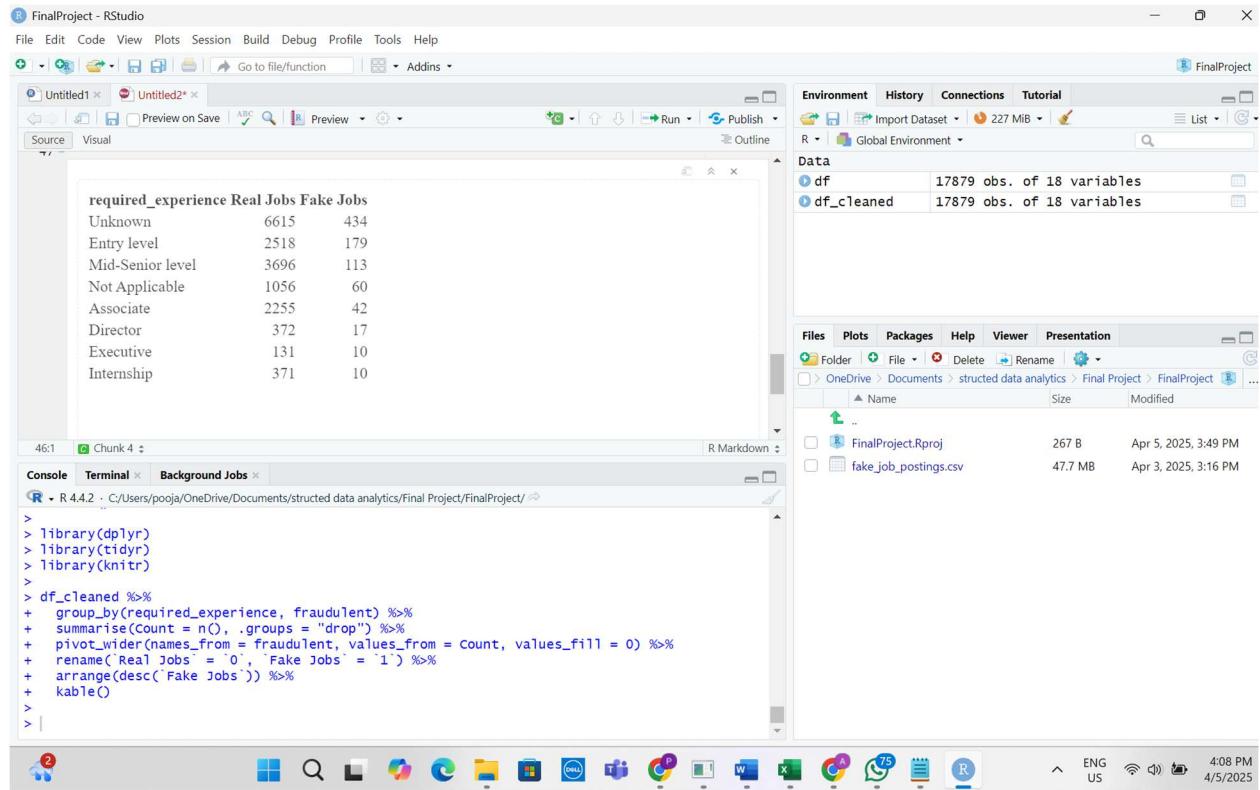


Figure 35: Pivot Table Based on Experience Required

Real vs Fake Jobs by Required Experience

Analyzing the job postings based on the level of experience required reveals significant differences in fraud distribution:

Key Observations:

- The “Unknown” experience level has the highest number of fake postings (434). This suggests that scammers frequently omit required experience to make the job listing appear accessible to all and reduce suspicion.
- Entry-level positions show the second-highest number of fake listings (179), which could be because scammers target people new to the job market who may be more vulnerable or less cautious.
- Mid-Senior level and Associate roles also show some fake listings, but at lower volumes, suggesting that scams are less common at more experienced levels.
- Positions like Director, Executive, and Internship have very few fake postings. This could be due to the smaller overall number of these listings and the increased scrutiny or specificity typically required for these roles.

Conclusion:

The data shows that job listings with ambiguous or missing experience requirements are more likely to be fraudulent. Entry-level roles are also frequently targeted, likely due to the broader applicant pool and higher response rates. This highlights required_experience as a valuable feature in fraud detection models.

3)

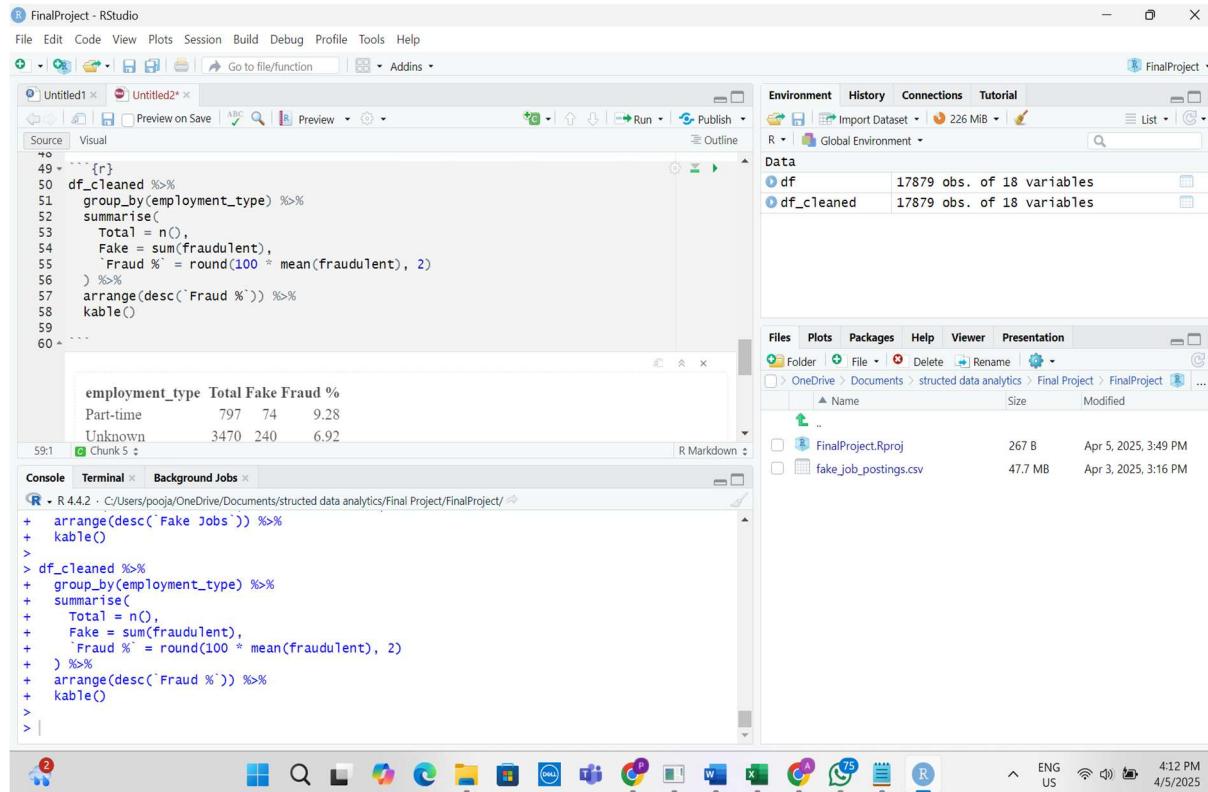


Figure 36: Pivot Table Based on Employment Type

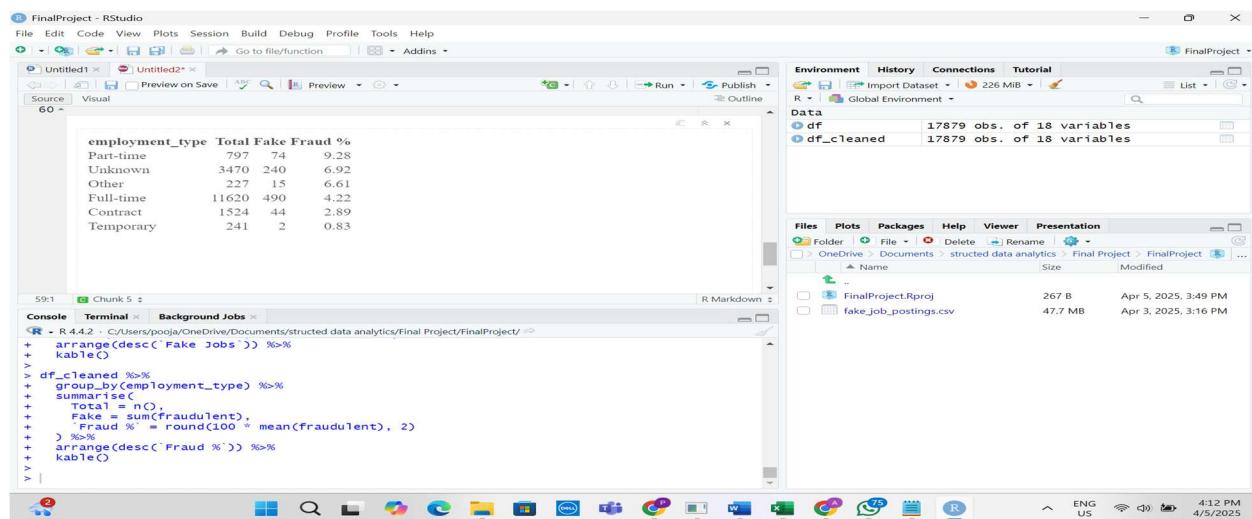


Figure 37: Results based on Pivot Table Based on Employment Type

Fraud Rate (%) by Employment Type

This pivot table highlights the percentage of job postings that are fraudulent across each employment type. Unlike raw counts, this view reveals how likely a job is to be fake depending on how it's described — and the results are very telling:

Key Findings:

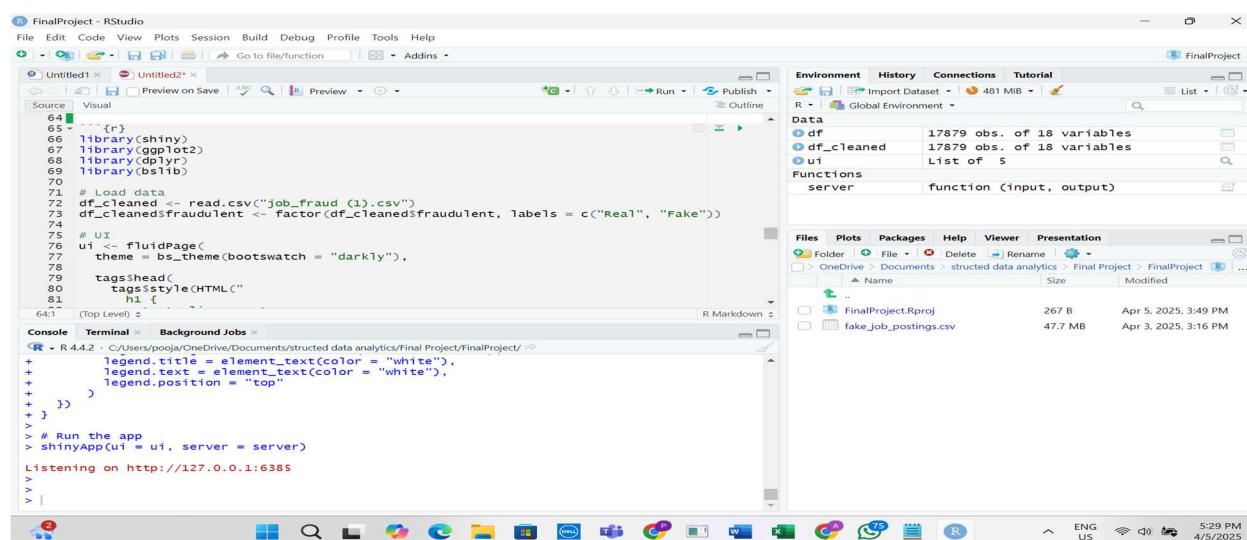
- Part-time jobs have the highest fraud rate at 9.28%, meaning nearly 1 in 10 part-time job postings are fake. This strongly suggests that scammers exploit the flexibility and popularity of part-time roles to attract more applicants.
- Unknown employment types also show a high fraud percentage (6.92%), reinforcing earlier insights that missing or vague job details are often associated with fake listings.
- Even though full-time jobs make up a majority of listings, the fraud rate is much lower at 4.2%. This shows that full-time jobs are generally more trustworthy — though scammers do still try to mimic them.
- Temporary and contract jobs have the lowest fraud rates, at 0.83% and 2.89% respectively. These roles might be less targeted because they usually involve formal hiring structures.

Conclusion:

This analysis makes it clear that fraud isn't just about volume — it's about proportion. While full-time jobs dominate the dataset, part-time and vague job types carry the most fraud risk per listing. These insights are crucial for building smarter fraud detection models that weigh fraud likelihood, not just counts.

EDA:

Dashboard in R or interactive plots:



The screenshot shows the RStudio interface with the following components:

- Source View:** Displays the R code for the dashboard. The code includes library imports (shiny, ggplot2, dplyr, bslib), data loading (df_cleaned from job_fraud_1.csv), creating a factor for fraudulence, defining a UI fluidPage with a darkly theme, and a server function that runs a shinyApp with the specified UI and server.
- Environment View:** Shows the global environment with objects like df (17879 obs. of 18 variables), df_cleaned (17879 obs. of 18 variables), and ui (List of 5).
- Files View:** Lists files in the project: FinalProject.Rproj (267 B) and fake_job_postings.csv (47.7 MB).
- Console View:** Shows the R session output, including the command to run the shinyApp and the resulting URL (http://127.0.0.1:6385).
- Top Bar:** Includes tabs for File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help, and Addins.

Figure 38: Dashboard Code

The screenshot shows the RStudio interface with the following components:

- Top Bar:** File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Source View:** Displays R code for a dashboard. The code includes CSS styling for h1 and .select-row classes, an HTML header for "Job Postings Fraud Detection Dashboard", and a div with class "select-row". It also contains R Markdown code for a legend and a shinyApp call.
- Console View:** Shows the R session output, including the R version (R 4.4.2), the current working directory (C:/Users/pooja/OneDrive/Documents/structured data analytics/Final Project/FinalProject), and the command to run the shinyApp.
- Environment View:** Shows the global environment with objects: df (17879 obs. of 18 variables), df_cleaned (17879 obs. of 18 variables), and ui (List of 5).
- Data View:** Shows the data frame df with 17879 observations and 18 variables.
- Files View:** Shows the project structure: OneDrive > Documents > structured data analytics > Final Project > FinalProject. It lists files: FinalProject.Rproj (267 B, Apr 5, 2025, 3:49 PM) and fake_job_postings.csv (47.7 MB, Apr 3, 2025, 3:16 PM).
- Bottom Bar:** Includes the taskbar with various application icons and system status indicators (ENG US, 5:30 PM, 4/5/2025).

Figure 39: Dashboard Code

The screenshot shows the RStudio interface with the following details:

- Title Bar:** FinalProject - RStudio
- File Menu:** File Edit Code View Plots Session Build Debug Profile Tools Help
- Source Editor:** Untitled1 (Visual mode) contains R code for creating a dashboard. The code includes sections for selecting rows, input fields for category and job type, and a bar plot output.
- Console:** Shows the R session output, including the command to run the shinyApp and the listening address (http://127.0.0.1:6385).
- Environment View:** Shows the global environment with objects df, df_cleaned, and ui.
- Files View:** Shows the project structure with files FinalProject.Rproj and fake_job_postings.csv.
- System Tray:** Shows battery level (ENG US), signal strength, and the date/time (4/5/2025, 5:30 PM).

Figure 40: Dashboard Code

The screenshot shows the RStudio interface with the following details:

- Title Bar:** FinalProject - RStudio
- File Menu:** File Edit Code View Plots Session Build Debug Profile Tools Help
- Source Editor:** Untitled1 (Visual mode) contains R code for creating a dashboard. The code includes sections for selecting rows, input fields for category and job type, and a bar plot output.
- Console:** Shows the R session output, including the command to run the shinyApp and the listening address (http://127.0.0.1:6385).
- Environment View:** Shows the global environment with objects df, df_cleaned, and ui.
- Files View:** Shows the project structure with files FinalProject.Rproj and fake_job_postings.csv.
- System Tray:** Shows battery level (ENG US), signal strength, and the date/time (4/5/2025, 5:30 PM).

Figure 41: Dashboard Code

The screenshot shows the RStudio interface with the following details:

- File Menu:** File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Source Editor:** Displays R code for a shiny application. The code includes data filtering, ggplot2 plotting (geom_bar), and a custom theme with white text on a dark background.
- Console:** Shows the command `shinyApp(ui = ui, server = server)` and the output "Listening on http://127.0.0.1:6385".
- Environment View:** Shows the global environment with objects like `df`, `df_cleaned`, and `ui`.
- Files View:** Shows files in the project directory: `FinalProject.Rproj` (267 B, modified Apr 5, 2025) and `fake_job_postings.csv` (47.7 MB, modified Apr 3, 2025).
- System Tray:** Shows battery level (ENG US), signal strength, and the date/time (5:30 PM, 4/5/2025).

Figure 42: Dashboard Code

The screenshot shows the RStudio interface with the following details:

- File Menu:** File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Source Editor:** Displays R code for a shiny application. The code includes data filtering, ggplot2 plotting (geom_bar), and a custom theme with white text on a dark background.
- Console:** Shows the command `shinyApp(ui = ui, server = server)` and the output "Listening on http://127.0.0.1:6385".
- Environment View:** Shows the global environment with objects like `df`, `df_cleaned`, and `ui`.
- Files View:** Shows files in the project directory: `FinalProject.Rproj` (267 B, modified Apr 5, 2025) and `fake_job_postings.csv` (47.7 MB, modified Apr 3, 2025).
- System Tray:** Shows battery level (ENG US), signal strength, and the date/time (5:31 PM, 4/5/2025).

Figure 43: Dashboard Code

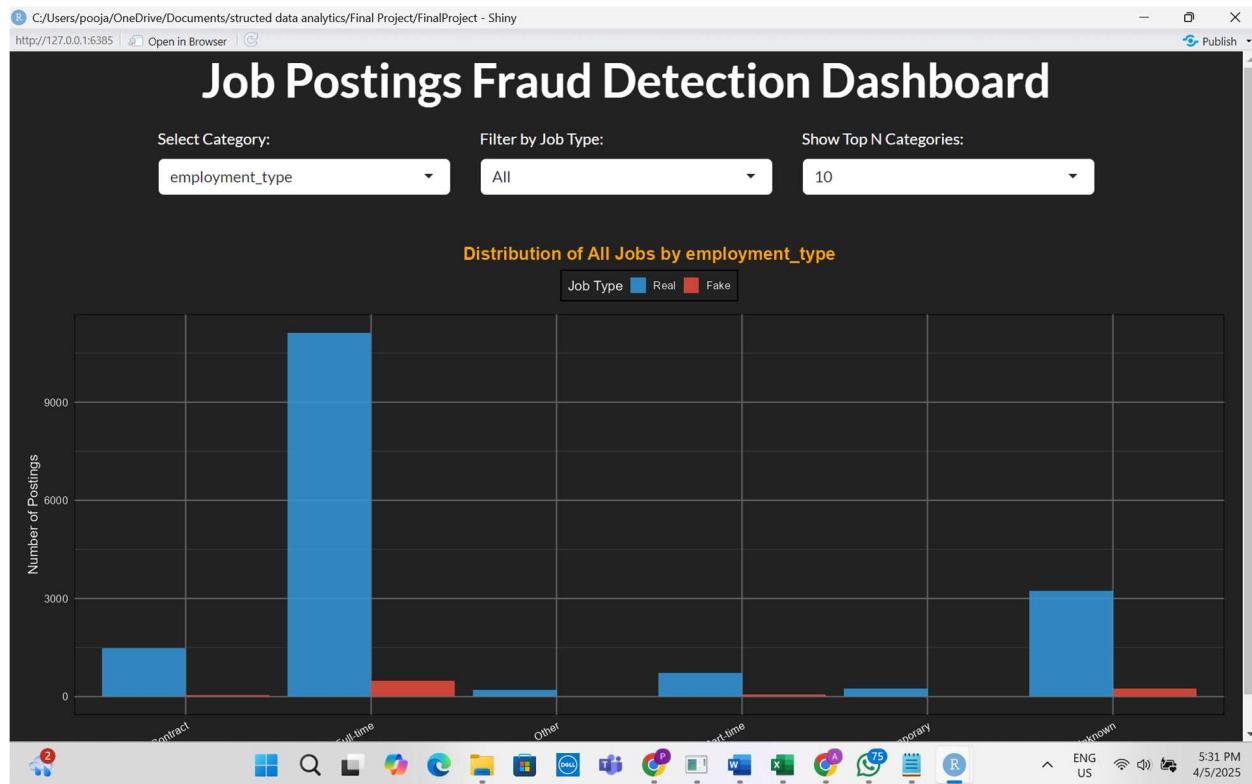


Figure 44: Dashboard



Figure 45: Dashboard

Interactive Dashboard Overview

This interactive dashboard was developed using R Shiny, based on a cleaned version of the *Fake Job Postings Dataset*. The data used in this dashboard was preprocessed to:

- Remove duplicate entries
- Drop null values in critical columns (like description)
- Impute or replace categorical nulls with “Unknown”
- Exclude irrelevant or redundant records to improve visualization clarity

The plot shown displays the distribution of real and fake job postings across various employment_type categories using a bar chart, with an option to filter by job type and category count.

Insights from the Employment Type Breakdown

- Full-Time jobs dominate the dataset, accounting for the largest share of both real and fake job postings.
 - While most full-time postings are real, a notable chunk of fake jobs also fall into this category, likely as scammers try to mimic professional job formats.
- The "Unknown" employment type shows a visible presence of fake postings.
 - This aligns with earlier findings from the ML model — vague or missing employment information is a red flag for fraud detection.
- Categories like Contract, Part-time, and Temporary are predominantly real, with very minimal fraud.
 - This may be due to their more specialized nature or clearer expectations, which scammers tend to avoid.
- The "Other" category contains very few listings overall but still includes some fakes — this might indicate inconsistent or non-standard job labeling, which again leans toward suspicious activity.

How This Dashboard Helps

- Filters like "Job Type" and "Top N Categories" enable users to interactively explore risk patterns by employment type, industry, and more.
- This can help HR professionals or platform moderators quickly spot suspicious categories, improving real-time fraud detection efforts.

Final Conclusion: Fake Job Postings Analysis

Summary of Key Findings

After performing a comprehensive analysis on the Fake Job Postings dataset from Kaggle, several patterns and risk factors that distinguish **fraudulent job postings** from **real ones** were identified. This investigation not only answered the core research questions but also demonstrated the effective use of Python and R to derive actionable insights.

1. What factors differentiate fake job postings from real ones?

- **Employment Type:** Fraudulent postings are disproportionately labeled as **Full-time**, likely to increase legitimacy. Real jobs show diversity across part-time, contract, and temporary categories.
- **Required Experience:** Fake jobs often avoid highly experienced roles. They tend to cluster around **Mid-Senior, Associate**, and especially **Unknown** experience levels to appeal to a broader group.
- **Industry & Function:** Fake listings are more common in **Information Technology, Internet, and Services**—industries with high online engagement. Rare job functions were found to be more prone to fraud.
- **Presence of Company Logo & Questions:** Real postings more often include logos, while frauds may include **screening questions** to appear legitimate.
- **Description Length & Language:** Fake job postings had a distinct word usage (e.g., "customer", "assist", "maintain") versus real ones ("team", "client", "experience"). Unusual length (too short or too long) also raised red flags.

2. Are there common attributes among fraudulent job listings?

Yes. Fraudulent job postings tend to share:

- **Missing or vague information** (like 'Unknown' in required experience or employment type)
- **Excessively generic or task-oriented language**
- **Absence of a company logo**
- **Unusual function labels** or rare job categories
- **Unrealistic or empty salary ranges**

- **Disproportionate frequency in certain industries and job types**

These attributes were quantified through EDA, violin plots, anomaly detection, and feature importance models.

3. Can data analysis techniques effectively detect fake job postings?

Absolutely. Your analysis used:

- **Outlier & anomaly detection (IQR, frequency-based)** to identify suspicious entries.
- **Correlation heatmaps** to identify patterns between features and the target variable (fraudulent).
- **Random Forest Classifier** with a well-structured pipeline, addressing class imbalance using `class_weight='balanced'`.

The model achieved:

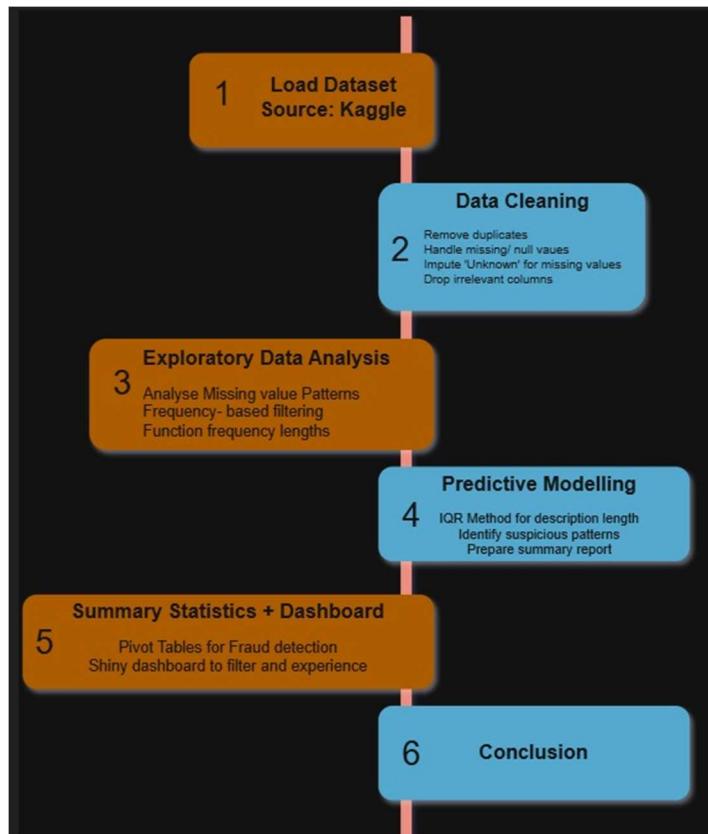
- **Accuracy:** 96%
- **F1 Score (fraud class):** 0.67 — A solid balance between false positives and false negatives.
- **Top Features:** function, industry, required_education, experience, and company_logo

These techniques together offer a **robust framework** to flag potentially fake postings before they reach job seekers.

Final Takeaways & Recommendations

- **Contextual awareness is critical:** Fraud patterns vary across industries and employment types. One-size-fits-all models are not effective.
- **Text analysis adds value:** Word clouds and NLP-based features can enhance fraud detection significantly.
- **Visualization improves clarity:** Through violin plots, pie charts, and heatmaps, patterns were made obvious and easier to interpret.
- **Interactive dashboards** (via R Shiny) empower stakeholders to explore data dynamically and take action faster.

Fraud Detection Pipeline for Job Postings



Why Python and R Were Used:

This project was completed using both **Python** and **R**, as required, but more importantly, their individual strengths were strategically used to handle different stages of the analysis efficiently.

Why Python Was Chosen for Data Cleaning, EDA, and Modeling

Python was used extensively for tasks that involved:

- **Data pre-processing:** Handling missing values, removing duplicates, and standardizing text entries was more intuitive using pandas.
- **Exploratory Data Analysis (EDA):** Python's seaborn and matplotlib libraries offer superior customization and interactivity for:
 - Detecting missing values (bar plots)

- Outlier detection (scatter plots, violin plots)
- Distribution plots for categorical and numerical features
- **Anomaly Detection:** Python made it easier to implement IQR-based filtering and frequency-based anomaly detection using standard statistical methods.
- **Text Analysis:** The WordCloud package in Python enabled quick visualization of word usage differences between real and fake jobs.
- **Machine Learning:** Python was ideal for model building and evaluation:
 - scikit-learn's Random Forest classifier with built-in support for handling imbalanced data
 - Evaluation metrics like confusion matrix, F1-score, precision/recall for fraud class

Python was chosen here because of its strong data manipulation capabilities, flexible visualization libraries, and robust machine learning frameworks that support end-to-end model development and interpretation.

Why R Was Chosen for Summary Statistics, Pivot Tables, and Interactivity

R was used in sections where **summary statistics and user-facing interpretation** were more important, such as:

- **Pivot Tables and Data Tables:**
 - R's dplyr and tidyverse libraries made it easier to generate clean summaries, group-wise counts, and frequency tables.
 - Quickly aggregating categorical variables to observe patterns (e.g., experience vs fraud status)
- **Interactive Visualizations:**
 - R's **Shiny dashboard** was used to create an interactive interface where users could:
 - Filter jobs by industry or experience
 - Visualize fraud patterns dynamically
 - Make the output more presentation-ready for stakeholders
- **Descriptive Analysis:**
 - R was useful for quick summaries, contingency tables, and exploratory insights to support the modeling narrative.

R was chosen in these parts due to its strong capabilities for statistical summaries, ease of working with categorical data, and most importantly, its ability to build lightweight, **interactive dashboards** via Shiny without additional setup.

References

1. **Dataset Source:**
 - o *Kaggle – Fake Job Postings Dataset.*
<https://www.kaggle.com/datasets/shivamb/real-or-fake-jobposting-prediction>
2. **Python Libraries:**
 - o *matplotlib: Visualization Library*
<https://matplotlib.org>
 - o *scikit-learn: Machine Learning in Python*
<https://scikit-learn.org>
3. **R Packages and Tools:**
 - o *Tidyverse – Data Science Framework*
<https://www.tidyverse.org>
 - o *Shiny – Web Application Framework for R*
<https://shiny.posit.co/r/getstarted/shiny-basics/lesson1/>
4. **Anomaly Detection Concepts:**
 - o *Wikipedia – Anomaly Detection*
https://en.wikipedia.org/wiki/Anomaly_detection