# RESUME PARSER

## ( END CAPSTONE PROJECT )

SUBMITTED BY: T.SASHIDHAR

COHORT: 20

INTEGRATED DATA SCIENCE

# TABLE OF CONTENTS

# 1. INTRODUCTION

## 1.1 Project Objective

The objective of this project is to develop an automated system for extracting, cleaning, storing, and analyzing resume data. The system leverages Python for data extraction and statistical analysis, SQL for storage and transaction management, Excel VBA for data cleaning, and Power BI for visualizations. The goal is to enhance the resume management and provide data-driveninsights that aid recruitment decision-making.

## 1.2 Project Deliverables

The deliverables for this project include:

- A Python script for extracting resume data from PDF and DOCX files.
- An Excel VBA macro for cleaning and standardizing the extracted resume data.
- SQL transactions, triggers, and stored procedures for data management.
- A Python class for statistical analysis of resume data.
- An XG Boost machine learning model to predict the primary skill from a resume.
- A Power BI dashboard that visualizes key resume metrics.

## 2.2 Tools and Technologies

- **Python Libraries**: PyPDF2, python-docx, Seaborn, NumPy, XGBoost, pandas, etc.
- **SQL**: SQL SMS for storage, transactions, triggers, and stored procedures.
- **Excel**: VBA for cleaning, Power Query for standardization.
- **Power BI**: Visualization and DAX calculations.

# 3. Data Extraction

## 3.1 Python Script for Extracting Resume Data

A Python script is written using PyPDF2 and python-docx to extract text from resumes. The extracted data is parsed into structured fields such as Name, Skills, Experience, and Education.



```
• EXTRACTING TEXT FROM THE DOCUMENT OR PDF FILE

[ ]   import os
      import re
      import PyPDF2
      import docx
      import openpyxl

      # FUNCTION TO EXTRACT TEXT FROM PDF FILES
      def extract_text_from_pdf(file_path):
          text = ""
          with open(file_path, 'rb') as file:
              reader = PyPDF2.PdfReader(file)
              for page_num in range(len(reader.pages)):
                  text += reader.pages[page_num].extract_text()
          return text

      # FUNCTION TO EXTRACT TEXT FROM DOCX FILES
      def extract_text_from_docx(file_path):
          doc = docx.Document(file_path)
          text = "\n".join([para.text for para in doc.paragraphs])
          return text
```
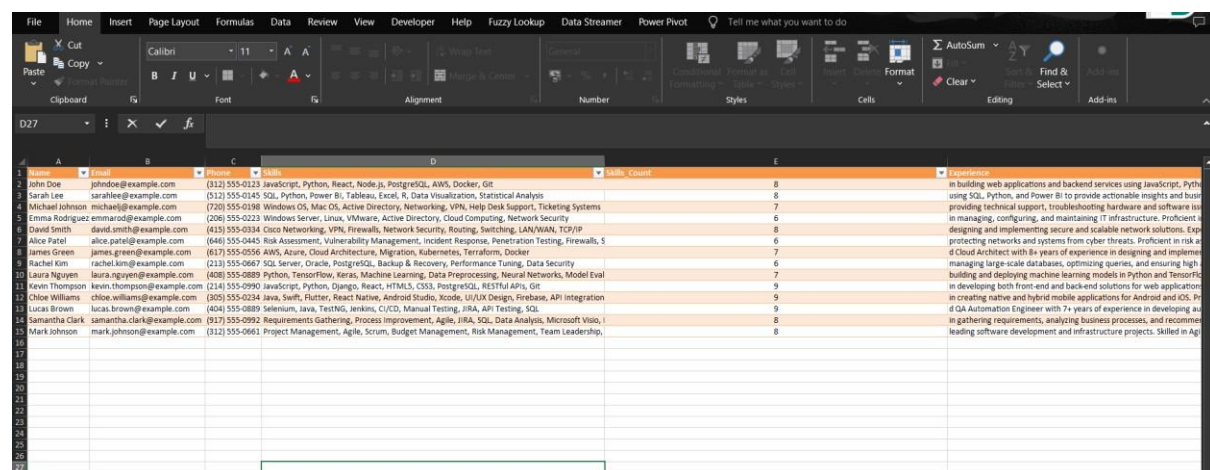
## 3.2 Mapping Extracted Data to Excel

The extracted data is then mapped to predefined columns (e.g., Name, Skills, Experience, Education) and stored in an Excel sheet.
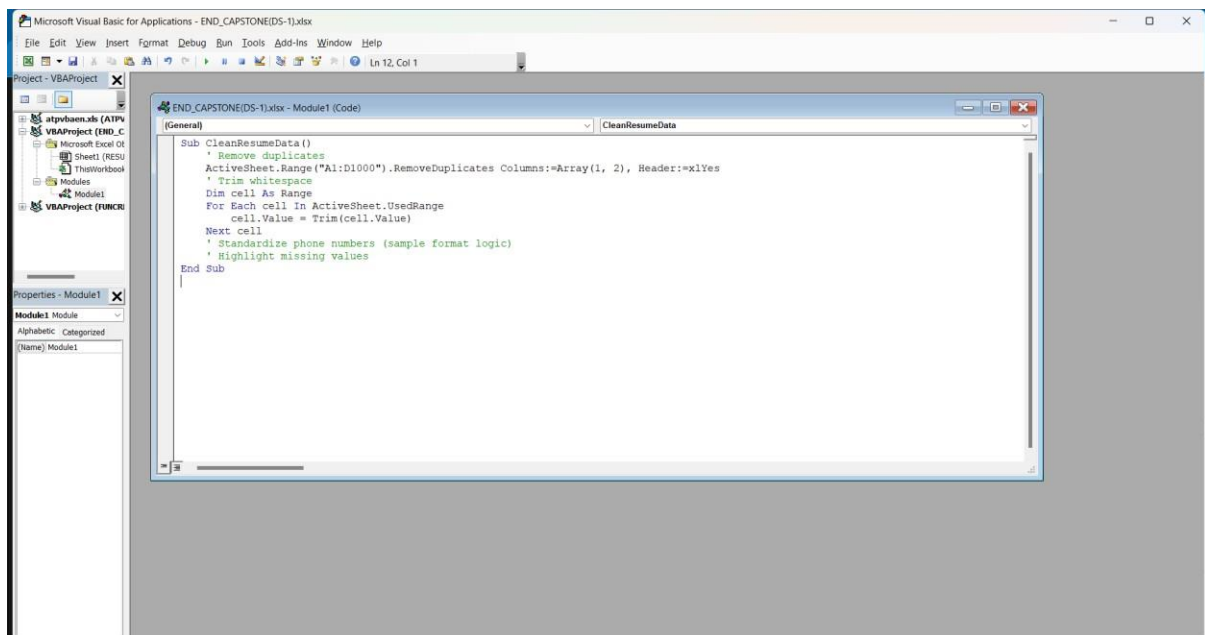
# 4. Data Cleaning

## 4.1 Excel VBA Macros for Data Cleaning

A VBA macro automates the cleaning process by:

- Removing duplicate entries.

- Trimming leading and trailing whitespaces from text fields.

- Standardizing phone numbers and email formats.

- Highlighting missing data.



## 4.2 Power Query for Data Standardization

Power Query is used to further standardize the resume data, ensuring consistent formats and filtering out incorrect entries.

# 5. Data Storage and Transaction Management

A SQL transaction is created to ensure data integrity. If an error occurs during insertion, the transaction is rolled back, ensuring that either all records are committed or none are.



## CREATE A SQL TRIGGER THAT AUTOMATICALLY LOGS THE INSERTION OF NEW RESUME DATA INTO A SEPARATE LOG TABLE

## STORED PROCEDURE THAT ACCEPTS MULTIPLE RESUME RECORDS AND INSERTS THEM INTO THE DATABASE IN A SINGLE CALL.





# 6. Statistical Analysis and Machine Learning

## 6.1 Python Resume Analyzer Class for Statistical Analysis

The Resume Analyzer class uses Seaborn and NumPy to perform statistical analysis on the resume data, such as analyzing the distribution of experience and the frequency of skills.

# Source Code:

## PYTHON CLASS RESUME ANALYZER THAT PERFORMS STATISTICAL ANALYSIS ON THE RESUME DATA. INCLUDE METHODS TO VISUALIZE THE DISTRIBUTION OF YEARS OF EXPERIENCE AND THE FREQUENCY OF DIFFERENT SKILLS

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import re

class ResumeAnalyzer:
    def __init__(self, data_path):
        """INITIALIZE THE RESUMEANALYZER WITH THE DATASET."""
        self.data = pd.read_excel(data_path)


    def experience_distribution(self):
        """PLOT THE DISTRIBUTION OF YEARS OF EXPERIENCE ACCURATELY."""
        def extract_experience(exp):
            if pd.isnull(exp):
                return np.nan
            # PATTERNS FOR VARIOUS FORMATS OF EXPERIENCE
            experience_match = re.search(r'\d+\+?\s*year', exp, re.IGNORECASE)
            if experience_match:
                exp = experience_match.group(0)
                if '-' in exp:
                    return float(exp.split('-')[0])
                # IF EXPERIENCE CONTAINS "LESS THAN A YEAR" OR SIMILAR
                elif 'less' in exp.lower() or 'fresher' in exp.lower():
                    return 0.0
                else:
                    # EXTRACT NUMBERS INCLUDIND DECIMAL PART
                    return float(re.findall(r'\d*\.?\d+', exp)[0])
            else:
                return np.nan

        # APPLY THE EXTRACTION FUNCTION TO THE 'EXPERIENCE' COLUMN
        self.data['Years_of_Experience'] = self.data['Experience'].apply(extract_experience)

        # DROPPING ROWS WITH NA VALUES
        experience_data = self.data['Years_of_Experience'].dropna()

        # PLOTTING THE DISTRIBUTION
        plt.figure(figsize=(10, 6))
        sns.histplot(experience_data, bins=10, kde=True, color='skyblue')
        plt.title('Distribution of Years of Experience')
        plt.xlabel('Years of Experience')
        plt.ylabel('Frequency')
        plt.show()

    def skills_frequency(self):
        """PLOT THE FREQUENCY OF DIFFERENT SKILLS."""
        all_skills = self.data['Skills'].dropna().str.split(',').sum()
        skill_counts = pd.Series(all_skills).value_counts()

        # PLOTTING THE TOP 10 MOST FREQUENT SKILLS
        plt.figure(figsize=(12, 8))
        sns.barplot(x=skill_counts.head(10).values, y=skill_counts.head(10).index, palette='viridis')
        plt.title('Top 10 Most Frequent Skills')
        plt.xlabel('Frequency')
        plt.ylabel('Skills')
        plt.show()


# OUTPUT
analyzer = ResumeAnalyzer("/content/drive/MyDrive/resume_data.xlsx")
analyzer.experience_distribution()
analyzer.skills_frequency()
```

# DISTRIBUTION OF YEARS OF EXPERIENCE AND THE FREQUENCY OF DIFFERENT SKILLS



```
<ipython-input-8-983d4e35f7fa>:58: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and set

  sns.barplot(x=skill_counts.head(10).values, y=skill_counts.head(10).index, palette='viridis')
```

```
sns.barplot(x=skill_counts.head(10).values, y=skill_counts.head(10).index, palette='viridis')
```



## 6.2 XGBoost Model for Skill Prediction

An XGBoost model is trained to predict the primary skill of a candidate based on resume content. The model is encapsulated in a Python class.

```
    Data loaded successfully.
                    Name                    Email          Phone  \
0          John Doe     johndoe@example.com  (312) 555-0123
1         Sarah Lee    sarahlee@example.com  (512) 555-0145
2   Michael Johnson    michaelj@example.com  (720) 555-0198
3   Emma Rodriguez      emmarod@example.com  (206) 555-0223
4       David Smith  david.smith@example.com  (415) 555-0334

                                             Skills  \
0   JavaScript, Python, React, Node.js, PostgreSQL...
1   SQL, Python, Power BI, Tableau, Excel, R, Data...
2   Windows OS, Mac OS, Active Directory, Networki...
3   Windows Server, Linux, VMware, Active Director...
4   Cisco Networking, VPN, Firewalls, Network Secu...

                                         Experience  \
0   in building web applications and backend servi...
1   using SQL, Python, and Power BI to provide act...
2   providing technical support, troubleshooting h...
3   in managing, configuring, and maintaining IT i...
4   designing and implementing secure and scalable...

                                          Education
0   Bachelor of Science in Computer Science – Univ...
1   Bachelor of Science in Data Science – Universi...
2   Associate of Science in Information Technology...
3   Bachelor of Science in Information Systems – U...
4   Bachelor of Science in Network Engineering – S...
Data preprocessed successfully.
Features extracted successfully.
Model trained successfully.
Model accuracy: 0.00
The predicted primary skill is: AWS, Azure, Cloud Architecture, Migration, Kubernetes, Terraform, Docker
/usr/local/lib/python3.10/dist-packages/xgboost/core.py:158: UserWarning: [16:52:14] WARNING: /workspace/src/learner.cc:740:
```

# 7. Data Visualization

## 7.1 Power BI Dashboard Overview

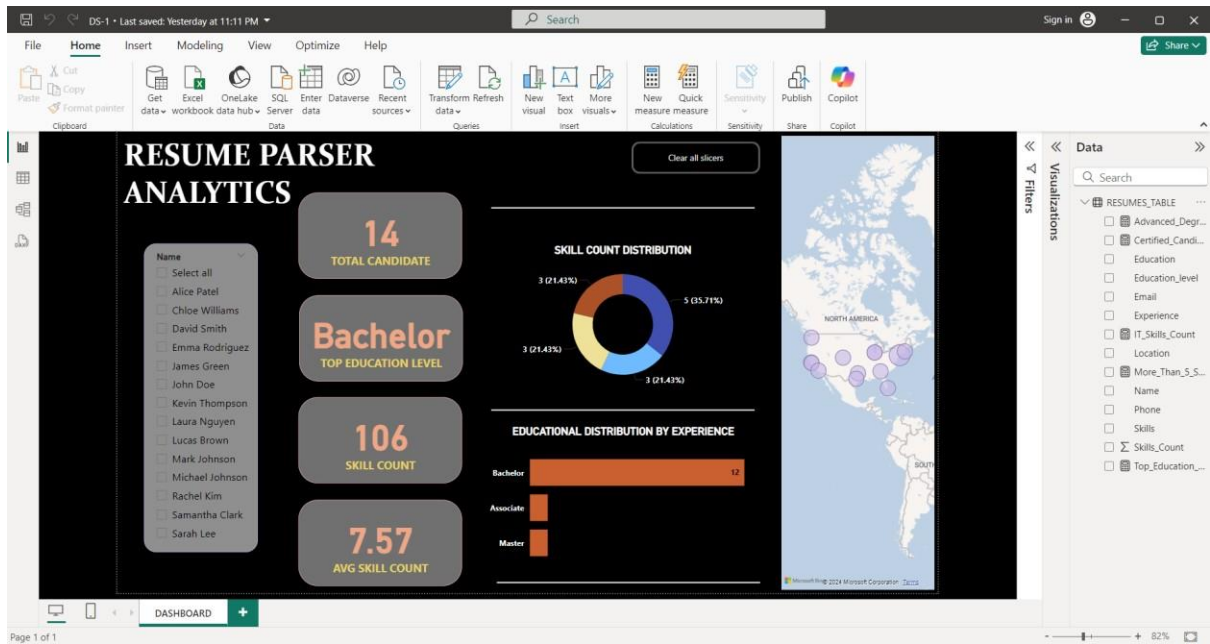A Power BI dashboard visualizes key metrics derived from the resume data, including:

- **Education Level Distribution**: Displays the highest education level of candidates.

- **Geographic Distribution**: Visualizes where candidates are located.

## 7.2 Key Metrics and Insights

- **Average Number of Skills per Candidate**: Helps identify candidates with diverse skills.

## 7.3 DAX Formulas for Calculations

Key metrics are calculated using DAX formulas, such as the average number of skills per candidate and the submission trends.

# 8. Conclusion

## 8.1 Summary of Achievements

The project successfully automates the extraction, cleaning, storage, and analysis of resume data using a combination of Python, SQL, Excel VBA, and Power BI.

**8.2 Key Insights Derived**

- **Skills**: Identified multi-skilled candidates.

- **Geographic Insights**: Helped understand where candidates are concentrated.

- **Trends**: Resume submission patterns inform recruitment strategies.

**8.3 Future Enhancements**

Future improvements could include expanding the machine learning model to predict other attributes such as job fit, and integrating external APIs to further enrich the resume data.