

Data Science Tools Project – Phase 1

1. Overview

This project focuses on automating the extraction, cleaning, and analysis of job listings from Wuzzuf, a leading online recruitment platform in Egypt. The main objective is to create a Python-based pipeline that collects data on job opportunities, extracts relevant details, cleans and organizes the data, and performs basic exploratory data analysis (EDA) to uncover trends and insights from the local job market.

2. Objectives

- Scrape job listings from Wuzzuf using Python.
- Extract essential job-related fields (title, company, location, skills, etc.).
- Clean and structure the data into a consistent format.
- Validate and categorize job fields to match Wuzzuf's known categories.
- Store the processed data in a database for reuse and scalability.
- Visualize trends to better understand job market demands.
- Deploy the analysis results as an interactive web app.

4. Data Collection

The scraping process is designed to:

- Iterate through multiple pages of Wuzzuf job listings.
- Extract fields such as:
 - Job Title
 - Company Name
 - Job Location
 - Employment Type (Full-time, Part-time)
 - Workplace Type (Remote, On-site, Hybrid)
 - Career Level
 - Years of Experience
 - Main Category
 - Skills
 - Job Link

To ensure completeness, the scraper continues to fetch pages until no new data is found. Each job is stored as a record in a Pandas DataFrame.

6. Data Storage

After processing, the final cleaned dataset is saved into a MongoDB database. This ensures that the data is not only persistent but also easily accessible for future use, especially if integrated into dashboards, recommendation engines, or APIs. MongoDB's flexible document-based structure makes it ideal for storing job listings, which may contain optional or nested fields like multiple skills or variable experience formats.

3. Tools and Technologies

The project utilizes the following tools and libraries:

- Python – for scripting and automation.
- Requests – for sending HTTP requests to fetch job listing pages.
- BeautifulSoup – for parsing HTML content and extracting data.
- Pandas – for data storage, cleaning, and manipulation.
- Matplotlib, Seaborn, Plotly – for visualization.
- Regex (re module) – for pattern matching and text cleaning.
- MongoDB – for storing the processed job listings in a NoSQL database.
- Streamlit – for deploying the analysis as an interactive web application.

5. Data Cleaning and Validation

- Empty or missing fields are labeled as "N/A" for clarity.
- Skills are cleaned and formatted as lists for easier analysis.
- The career level, employment type, and category fields are cross-checked and mapped to Wuzzuf's standardized categories.
- Jobs with invalid or unrecognized categories are either fixed using a custom mapping or excluded from final analysis.

7.Exploratory Data Analysis (EDA)

- To understand market demands and job trends, EDA was conducted using visual tools:
- Bar charts show the distribution of top-requested skills.
- Pie charts illustrate the proportions of career levels across listings.
- Histograms highlight how many years of experience are typically required.
- Category frequency plots reveal the most popular job fields.
- These visualizations provide valuable insights for job seekers, recruiters, and analysts who want to track trends in the Egyptian tech job market.

8. Challenges and Solutions

✓ Challenges:

- HTML inconsistencies between job postings.
- Multiple or missing categories and career levels.
- Duplicate job listings across pages.
- Free-form skill tags requiring standardization.

🔧 Solutions:

- Used robust HTML parsing and fallback mechanisms to extract fields.
- Built a filtering system to validate and correct job categories.
- Applied regular expressions and logic to extract structured information from unstructured text.

9. Bonus Task: Interactive Web App Deployment

As a bonus step, the processed data and analytical results were deployed using Streamlit, a Python-based framework for building interactive web applications. This allows users to:

- Explore top skills and job categories dynamically.
- Filter listings based on experience level, skills, or location.
- Visualize trends directly from the browser without any coding experience.

This deployment adds a user-friendly interface to the backend analysis, enabling real-time exploration and sharing of insights.

10. Conclusion

This project successfully demonstrates how to build an end-to-end job scraping and analysis pipeline. It scrapes real job data, transforms it into structured and meaningful datasets, stores it in a database, and uncovers useful patterns through visualization. The deployment of results in an interactive app further extends the project's reach and usability.