

Wuzzuf Job Data Cleaning Project Documentation

1. Executive Summary

This project implements an automated data pipeline to analyze the tech job market in the MENA region. The workflow consists of three main phases:

1. **Data Collection:** A Python automation tool scrapes real-time job data from Wuzzuf for 10 specific technical roles.
2. **Data Cleaning:** An ETL (Extract, Transform, Load) process standardizes the raw data, making it suitable for analysis.
3. **Data Analysis:** Exploratory analysis to visualize trends, geographical distribution, and role demand.

Input: Dynamic web content from Wuzzuf. Final Output: Cleaned_jobs.csv (Structured Dataset) & Analytical Insights.

2. Phase 1: Data Collection (Web Scraping)

2.1 Tools & Libraries

- **Python 3:** Core programming language.
- **BeautifulSoup (bs4):** Parses HTML content to extract job titles, companies, and locations.
- **Requests:** Handles HTTP requests to retrieve page content.
- **Urllib:** Manages dynamic URL encoding for search queries.

2.2 Implementation Logic

- **Target Scope:** The script iterates through 10 key roles: *Data Engineer, Data Analyst, Data Scientist, AI Engineer, Backend Developer, Frontend Developer, Software Engineer, Network Engineer, DevOps Engineer, Full Stack Developer*.
- **Pagination:** A while loop automatically navigates through all available search result pages until no jobs remain.
- **Extraction:** Specific CSS selectors are used to pull data points from job cards, with error handling for missing fields.

3. Phase 2: Data Cleaning & Preprocessing

Wuzzuf Job Data Cleaning Project Documentation

3.1 Tools & Libraries

- **Pandas:** For data manipulation and CSV I/O.
- **NumPy:** Handles numerical arrays and missing values (NaN).
- **Re (Regex):** Performs pattern matching for text cleaning.
- **String:** Provides string constants for punctuation removal.

3.2 Transformation Steps (ETL)

The raw scraped data undergoes the following cleaning process:

1. Column Selection (Dropping Noise):

- **Dropped:** Area (Too granular/many missing values), Publish Time (Relative time is ambiguous), Job Link (Not needed for aggregate analysis).
- **Kept:** Core attributes like Company, Location, and Job Title.

2. Standardization (Renaming): Columns were renamed to Python snake_case for consistency:

- Searched Job → Job
- Job Title → Job_title
- Job Type → Job_type
- Work Place → Work_place

3. Text Cleaning:

- A custom function removes punctuation and special characters from Job_title.
- **Goal:** Ensures "Data Engineer" and "Data Engineer - Backend" are treated as textually similar.

4. Handling Missing Values:

- **Work_place:** Missing values are filled with "Undefined" to ensure data completeness.

Wuzzuf Job Data Cleaning Project Documentation

4. Phase 3: Exploratory Data Analysis (EDA)

4.1 Tools & Libraries

- **Matplotlib:** For creating static, animated, and interactive visualizations.
- **Seaborn:** For making statistical graphics and heatmaps.

4.2 Key Insights

Based on the analysis of the `Cleaned_jobs.csv` dataset, the following trends were identified regarding the relationship between Country and Job Role:

1. **Dominant Market:** Egypt shows the highest hiring activity across all technical job categories, serving as the primary hub for the platform.
2. **Regional Tech Hubs:** Saudi Arabia and the United Arab Emirates follow as strong markets for technical talent.
3. **Most In-Demand Role:** "Software Engineer" is the most frequently posted job title across all analyzed countries.
4. **Specialized Demand:** Data-related roles (Data Engineer, Data Scientist) show strong, concentrated demand in specific regions.
5. **Low Activity Regions:** Countries such as Belgium, Georgia, Indonesia, Libya, Qatar, and India appeared in the search results but with negligible volume (0–2 posts), indicating they are outliers on this specific platform.

5. Data Dictionary

This section outlines the schema transformation from the raw scraped data to the final clean dataset.

5.1 Raw Data (Input: `jobs_wuzzuf.csv`)

Column Name Description

Searched Job The keyword used in the search query (e.g., "Data Engineer").

Job Title Raw title from the ad (often contained punctuation/emojis).

Wuzzuf Job Data Cleaning Project Documentation

Company	Name of the hiring company.
Location	Raw location strings parsed into Country, City, and Area.
Publish Time	Relative string (e.g., "7 days ago").
Job Link	Direct URL to the post.
Job Type	Employment status (Full-time, Freelance).
Work Place	Modality (Remote, On-site, Hybrid).

5.2 Cleaned Data (Output: Cleaned_jobs.csv)

Column Name Description

Job	Standardized category (e.g., Data Engineer).
Job_title	Cleaned text (Punctuation removed).
Company	Preserved company name.
Country	Preserved country.
City	Preserved city.
Job_type	Employment status.
Work_place	Modality (NaNs filled with "Undefined").
Scraping date	Date the data was collected (Primary time reference).

6. Conclusion

This project successfully demonstrates an end-to-end data pipeline. By utilizing Python and BeautifulSoup, unstructured web data was transformed into a structured dataset. Subsequent analysis revealed that while "Software Engineer" remains the universal staple of the tech industry, Egypt, Saudi Arabia, and the UAE are the definitive centers of gravity for tech employment on the Wuzzuf platform.