# M.F.P Assignment – Report

========================================================================

**Author:** Aditya Kayasth

**Date:** 30-Nov-25

========================================================================

## Task:

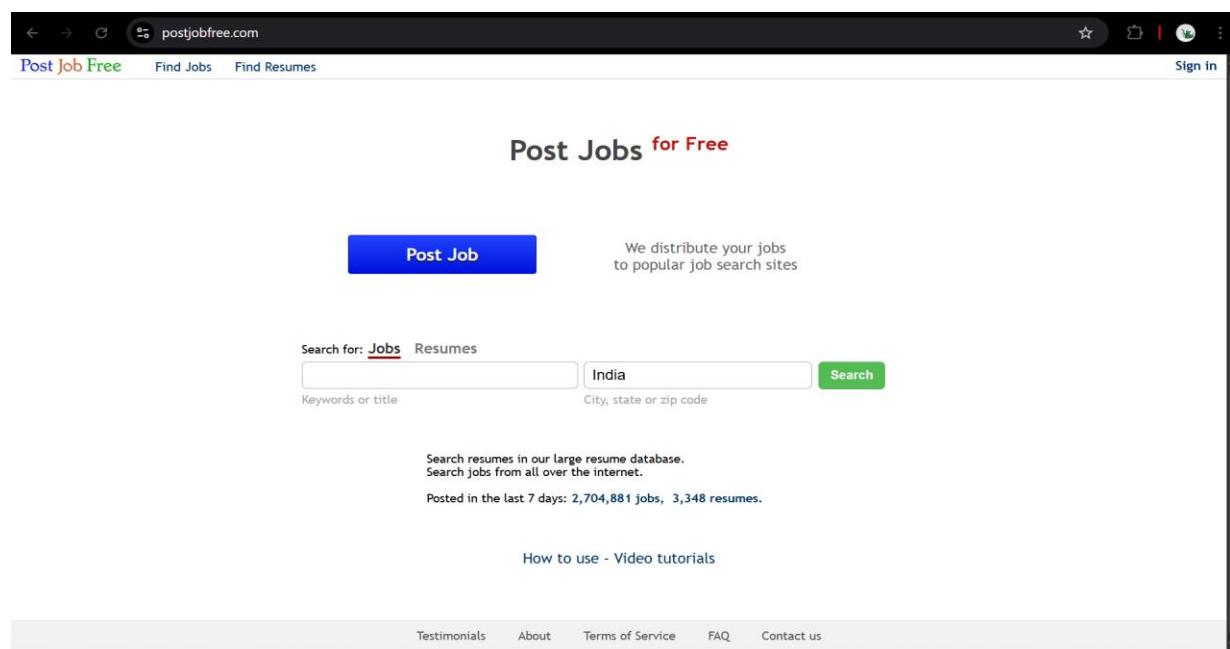Create Candidate Search API (Skill + Experience Filter)

## Objective:

Build an API endpoint that accepts skill name and number of years of experience as input parameters and returns a list of matching candidate profiles. The candidate data should be sourced from Naukri.com (via scraping and available APIs) and any other reliable publicly available sources (1 source) that allow candidate data retrieval.

========================================================================

## Data source used:

### 1) PostJobFree.com

**Description:** It is an open database where candidates upload their resumes specifically to be found by recruiters. They do not use aggressive anti-bot technology. They allow scraping to read their pages using the normal and advanced queries via request module.

**Query Structure:**

1) Normal: https://www.postjobfree.com/resumes?q={skill}&l={location}
2) Advanced: https://www.postjobfree.com/advanced-resume-search?q=title:({job_title})+({must_have_skills})({excluded_keywords})&l={location}&radius={proximity}&r={limit}

**2) Naukri.com (Primary Target):**



**Note:** Direct scraping of Naukri is strictly blocked as it violates their terms. They monetize this data via "Resdex" (priced at ~₹4,000 for 100 views), making direct unauthorized access difficult and unethical.

=============================================================================

## Challenges:

### 1) TLS/SSL Fingerprinting:

During testing, Google X-Ray requests were frequently blocked, even when rotating User-Agents. Standard Python requests libraries leak their identity via the TLS Handshake (JA3 Fingerprint). Google detects that the request is coming from a script (OpenSSL) rather than a real browser.

### Solution:

Implementing advanced libraries such as **curl_cffi**, **Selenium**, or **Scrapy** allows the scraper to impersonate real browser TLS signatures, significantly reducing detection rates.

### 2) Data Consistency:

Data returned from public sources like PostJobFree can be inconsistent. Fields like "Years of Experience" vary in format, and sensitive contact information is often redacted in public search views.

### Example:

### API REQUEST:

```
{
    "skill": "java",
    "experience": 3
}
```

### API RESPONSE (Truncated):

```
{
    "status": "success",
    "count": 10,
    "candidates": [
        {
            "source": "PostJobFree",
            "name": "Candidate (Hidden)",
            "current_job_title": "Academic-Driven Assistant Professor with ML & CSE Focus",
```

      "skills": ["java"],

      "experience_years": "Check Profile",

      "location": "India",

      "resume_url": "https://www.postjobfree.com/resume/aef5wo..."

    }

  ]

}

**RESUME URL RESULTS:**



**Solution:** I implemented a Regex-based parser to standardize experience data into integers for filtering. To retrieve complete profiles, a secondary scraper that could be deployed to visit individual resume URLs. Integrating a Large Language Model (LLM) would further enhance the extraction of unstructured data from these full resume texts.

===========================================================================

## Improvements planned:

1) **Scraping Efficiency:**

I plan to fix the blocking issues by upgrading from standard Python requests to curl_cffi. This tool impersonates real browsers, which stops Google and Naukri from blocking the connection. I would also add Proxy Rotation to switch IP addresses automatically, so I can search thousands of times without getting banned.

2) **Data Consistency:**

I want to fetch better data. Currently, the API retrieves search summaries and no contact details about the candidate can be found at all.  I plan to build a scraper that visits every resume link to grab full work history and contact info. To keep this fast, I would use Multithreading to download multiple profiles at the same time, and maybe use an LLM (AI) to clean up the messy text automatically.

========================================================================

**Gitlab Link:** https://gitlab.com/internship-task1/MFP-Task-API-V1/-/tree/main

========================================================================