# *Indeed Job Scraper*

## 1. Introduction

Searching for jobs online has become an essential step for students, professionals, and recruiters. However, manually browsing job boards like **Indeed** is repetitive and time-consuming. Job seekers may need to analyze trends such as which companies are hiring more, which cities have more opportunities, or which skills are in demand.

The **Indeed Job Scraper with Simple Analytics** is a Python-based mini project designed to automate this process. It extracts job details from Indeed and generates **data-driven insights** using charts and word clouds.

## 2. Problem Statement

- Job seekers manually check job boards every day, wasting time and effort.

- Recruiters and students lack quick tools to identify **hiring companies, locations, and skill trends**.

- A simple scraper can solve this problem but must also provide **unique features** to be useful.

## 3. Objectives

The objectives of this project are:

1. Automate job search data collection from **Indeed**.

2. Extract **job title, company, location, date posted, description, and job link**.

3. Save the results in **CSV/Excel format** for further use.

4. Provide **extra features** such as:

- ○ Visualization of **top hiring companies**.

- ○ **Jobs by location** analysis.

- ○ A **Word Cloud** to highlight trending skills from job descriptions.

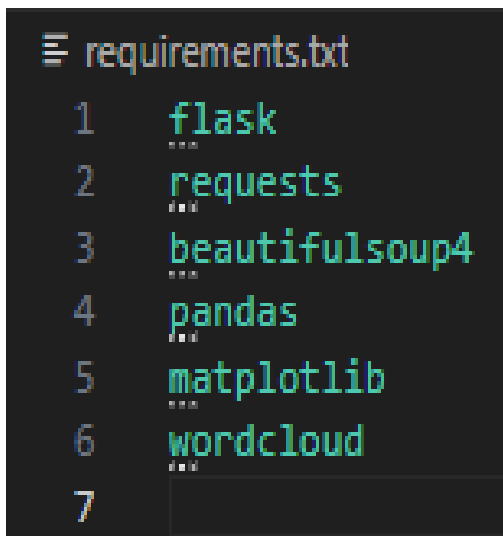- ○ **Simple filtering** of jobs by company or location.

**4. Tools and Technologies**

**4.1 Programming Language**

- **Python 3.9+**

**4.2 Libraries Used**

- **requests** → Fetch HTML pages from Indeed.

- **beautifulsoup4** → Parse and extract job details from HTML.

- **pandas** → Store, clean, and export results into CSV/Excel.

- **matplotlib** → Generate bar charts for analytics.

- **wordcloud** → Create word clouds of skills/keywords from job descriptions.

```
requirements.txt
1    flask
2    requests
3    beautifulsoup4
4    pandas
5    matplotlib
6    wordcloud
7
```

# 5. System Design

## 5.1 Workflow

```
User Input (Keyword + Location)
        ↓
Requests → BeautifulSoup
        ↓
Extract job details (Title, Company, Location, Date, Description, URL)
        ↓
Save in Pandas DataFrame
        ↓
Export as CSV/Excel
        ↓
Analytics
    - Top companies (bar chart)
    - Jobs by location (bar chart)
    - Skill cloud (word cloud)
        ↓
Optional: Filter jobs by company/location
```

## 5.2 Data Model (columns saved in CSV)

- Job Title
- Company
- Location
- Date Posted
- Job Description
- Job URL

# 6. Implementation

## 6.1 Data Collection

- Construct the Indeed search URL using keyword and location.
- Use requests to fetch job listing pages.
- Parse HTML with BeautifulSoup to extract details.

- Loop through 2–3 pages using pagination**.**

## 6.2 Data Storage

- Store scraped job data in a pandas DataFrame.
- Export results into:
    - jobs.csv (CSV format)
    - jobs.xlsx (Excel format, optional**)**

## 6.3 Data Analytics (Unique Features)

1. Top Hiring Companies
    - Count job postings by company.
    - Plot a bar chart of top 5 companies.
2. Jobs by Location
    - Count job postings by city/region.
    - Plot a bar chart of top 10 locations.
3. Skill Word Cloud
    - Combine all job descriptions.
    - Generate a word cloud showing frequent terms (e.g., *Python, SQL, Excel*).
4. Simple Filtering
    - Allow users to filter the CSV results by company or location.

---

## 7. Sample Outputs

## 7.1 CSV File (jobs.csv)

| Job Title | Company | Location | Date Posted | Description | Job URL |
|---|---|---|---|---|---|
| Data Analyst | TCS | Chennai | 2 days ago | Analyze data... | … |
| Python Developer | Infosys | Bangalore | 5 days ago | Develop Python apps... | … |

## 7.2 Visualization (Matplotlib Charts)

- Top 5 Hiring Companies – Bar chart showing which companies have the most postings.
- Jobs by Location – Bar chart showing job counts in different cities.

## 7.3 Word Cloud (Skills)

- Word cloud image generated from job descriptions.
- Highlights common skills like Python, SQL, Excel, Machine Learning.

---

# 8. Results

- Successfully scraped job data from Indeed for a chosen keyword and location.
- Stored structured job data in CSV format.
- Generated insights:
    - Top companies hiring.
    - Job distribution by location.
    - Trending skills in demand.
- Added filtering for customized job search.

---

# 9. Applications

- Students: Identify trending skills for placements.
- Job Seekers: Save time by analyzing multiple job postings at once.
- Recruiters: Observe competitor hiring trends.
- Researchers: Study market demand across industries.

---

# 10. Limitations

- Some jobs may not have salary or full descriptions available.
- Website HTML structure may change (scraper must be updated).
- Heavy scraping may trigger rate limits

---

# 11. Future Enhancements

- Add salary extraction for salary trend analysis.
- Use Selenium for dynamic scraping.
- Automate scraping daily using schedule/cron jobs.
- Build a Streamlit web app for interactive job search.
- Compare multiple roles (e.g., Data Analyst vs Software Engineer).

TEAM-6
MAHALAKSHMI S

## 12. Conclusion

The Indeed Job Scraper with Simple Analytics is a practical mini project that goes beyond plain web scraping. By combining data extraction with visualization and word clouds, it transforms job postings into actionable insights. This project is simple enough to implement yet unique enough to stand out.

# Cybernaut mini project-1

# Appendix

# Scarper.py

```python
scraper.py > ...
1   import requests
2   from bs4 import BeautifulSoup
3   import pandas as pd
4   import time, random, urllib.parse, os
5   import matplotlib.pyplot as plt
6   from wordcloud import WordCloud, STOPWORDS
7
8   BASE_URL = "https://www.indeed.com/jobs"
9
10  HEADERS = {
11      "User-Agent": (
12          "Mozilla/5.0 (Windows NT 10.0; Win64; x64) "
13          "AppleWebKit/537.36 (KHTML, like Gecko) "
14          "Chrome/115.0 Safari/537.36"
15      ),
16      "Accept-Language": "en-US,en;q=0.9",
17      "Referer": "https://www.google.com/",
18      "Connection": "keep-alive",
19  }
20
21  if not os.path.exists("static"):
22      os.makedirs("static")
23
24  def scrape_jobs(keyword, location, pages=2):
25      """Scrape jobs from Indeed. Fallback to sample dataset if blocked."""
26      all_jobs = []
27      for i in range(pages):
28          params = {"q": keyword, "l": location, "start": i * 10}
29          r = requests.get(BASE_URL, params=params, headers=HEADERS, timeout=15)
30
31          if r.status_code == 403 or "job_seen_beacon" not in r.text:
32              print("⚠ Blocked by Indeed or no results. Loading sample dataset.")
33              return pd.read_csv("sample_jobs.csv", quotechar='"', on_bad_lines="skip")
34
35          soup = BeautifulSoup(r.text, "html.parser")
36          cards = soup.select("div.job_seen_beacon")
37          if not cards:
38              print("⚠ No job cards found. Using sample dataset.")
39              return pd.read_csv("sample_jobs.csv", quotechar='"', on_bad_lines="skip")
40
41          for card in cards:
42              title_tag = card.select_one("h2.jobTitle") or card.select_one("a.jobtitle")
43              company_tag = card.select_one(".companyName") or card.select_one(".company")
44              loc_tag = card.select_one(".companyLocation")
45              date_tag = card.select_one(".date")
46              desc_tag = card.select_one(".job-snippet") or card.select_one(".summary")
47              a_tag = card.select_one("a")
48
49              all_jobs.append([
```
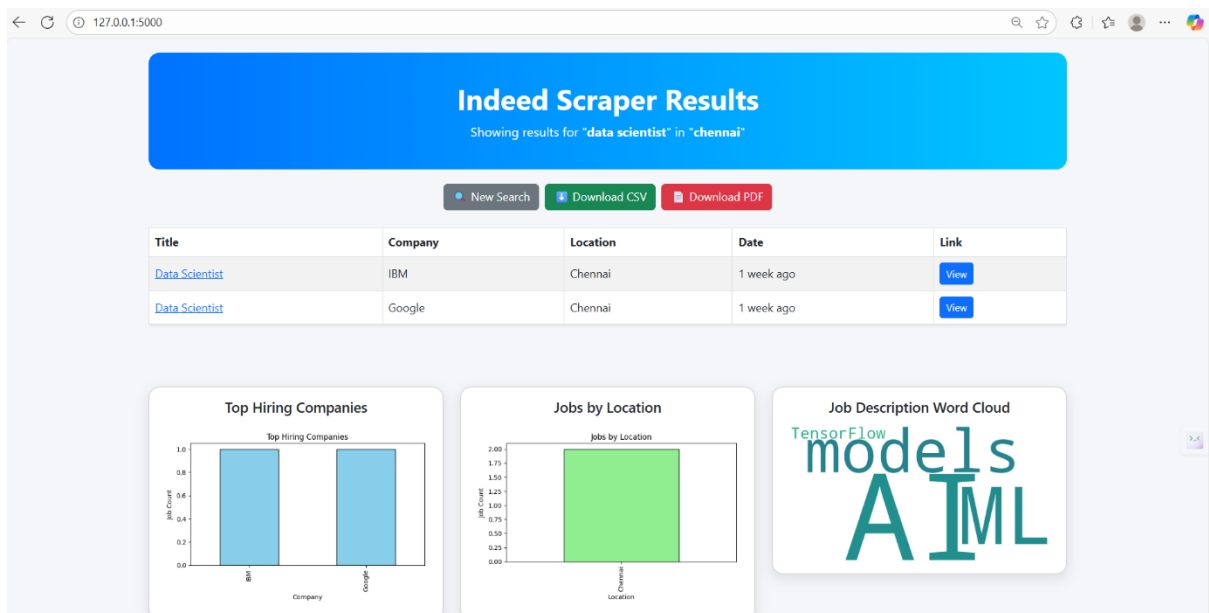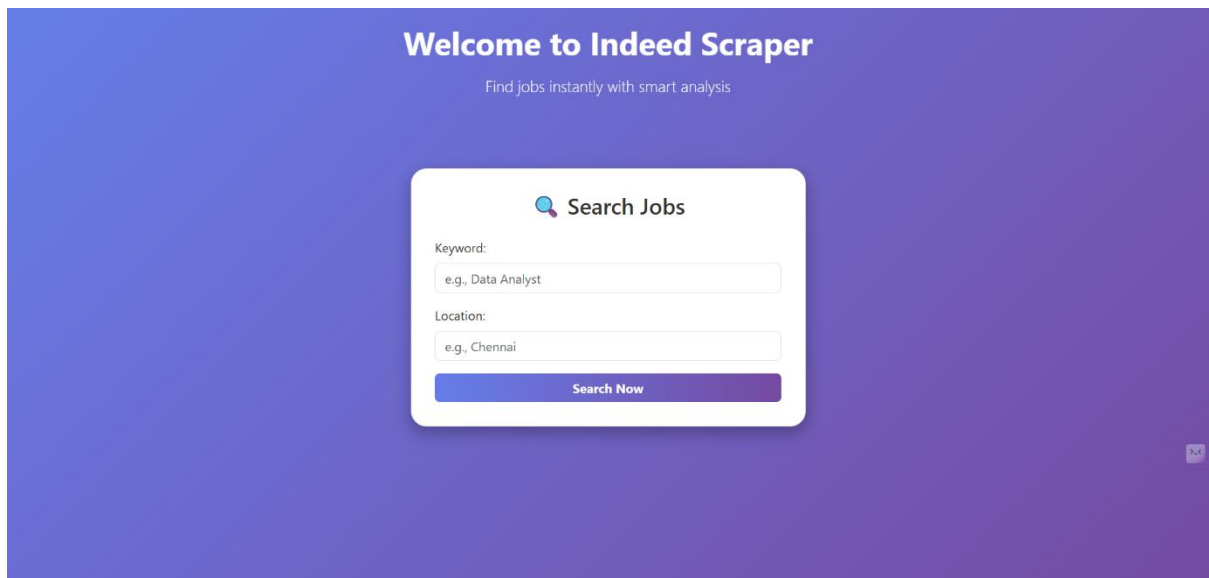
# App.py

```
app.py > download_pdf
  1   from flask import Flask, render_template, request, send_file
  2   from scraper import scrape_jobs, generate_charts
  3   import os
  4   import pandas as pd
  5   from reportlab.lib.pagesizes import letter
  6   from reportlab.pdfgen import canvas
  7   app = Flask(__name__)
  8   @app.route("/", methods=["GET", "POST"])
  9   def index():
 10       if request.method == "POST":
 11           keyword = request.form.get("keyword").lower()
 12           location = request.form.get("location").lower()
 13           df = scrape_jobs(keyword, location, pages=2)
 14           if not df.empty:
 15               df = df[
 16                   df["Title"].str.lower().str.contains(keyword, na=False) &
 17                   df["Location"].str.lower().str.contains(location, na=False)
 18               ]
 19           if df.empty:
 20               return render_template(
 21                   "results.html",
 22                   jobs=[],
 23                   keyword=keyword,
 24                   location=location,
 25                   no_results=True,
 26                   top_companies_exists=False,
 27                   jobs_by_location_exists=False,
 28                   wordcloud_exists=False
 29               )
 30           df.to_csv("jobs.csv", index=False)
 31           generate_charts(df)
 32           jobs = df.to_dict(orient="records")
 33           context = {
 34               "jobs": jobs,
 35               "keyword": keyword,
 36               "location": location,
 37               "no_results": False,
 38               "top_companies_exists": os.path.exists("static/top_companies.png"),
 39               "jobs_by_location_exists": os.path.exists("static/jobs_by_location.png"),
 40               "wordcloud_exists": os.path.exists("static/wordcloud.png"),
 41           }
 42           return render_template("results.html", **context)
 43       return render_template("index.html")
 44   @app.route("/download_csv")
 45   def download_csv():
 46       if os.path.exists("jobs.csv"):
 47           return send_file("jobs.csv", as_attachment=True)
 48       return "⚠ No jobs file found!"
 49   @app.route("/download_pdf")
 50   def download_pdf():
 51       if not os.path.exists("jobs.csv"):
 52           return "⚠ No jobs file found!"
 53       file_path = "jobs.pdf"
 54       df = pd.read_csv("jobs.csv")
 55       c = canvas.Canvas(file_path, pagesize=letter)
 56       width, height = letter
 57       c.setFont("Helvetica-Bold", 16)
 58       c.drawString(200, height - 50, "Job Listings Report")
 59       c.setFont("Helvetica", 10)
 60       y = height - 80
 61       for _, row in df.iterrows():
 62           text = f"{row['Title']} | {row['Company']} | {row['Location']} | {row['Date']}"
 63           c.drawString(50, y, text)
 64           y -= 15
 65           if y < 50:
 66               c.showPage()
 67               c.setFont("Helvetica", 10)
```

# Cybernaut mini project-1

## Output:





TEAM-6
MAHALAKSHMI S

# Cybernaut mini project-1

TEAM-6
MAHALAKSHMI S