# INDEED SCRAPER

PYTHON BASED

# TEAM 6

# TEAM MEMBERS

PADMAJA S.P (TEAM LEAD)
SUTHIKSHA .S
KAVISH SHARMA
SANDHIYA .S
RAGHUNANDHAN .K
SHALINI .K

## PROJECT INTRODUCTION

The objective of this project is to develop a web scraping tool that can automatically extract job listing data from the Indeed website. The tool will collect a wide range of information, including job titles, company names, locations, salaries (if available), job descriptions, and posting dates. This data can be used for analysis, trend identification, and various applications such as job market analysis, salary comparisons, and job trend forecasting

JOB TITLE EXTRACTION
COMPANY NAME EXTRACTION
LOCATION
SALARY INFORMATION
JOB DESCRIPTION

# TOOLS & TECHNOLOGIES USED

**Python**: The programming language used to implement the scraper logic.

**Selenium**: A web automation tool used to interact with web pages and extract data

**dynamically.Chromedriver**: A browser driver that allows Selenium to control Google

**Chrome.logging**: A Python library used for tracking events and debugging by recording
**runtime messages.os**: A Python library to interact with the operating system for file paths and
**environment variables.re**: A Python library for handling regular expressions, used for string
**pattern matching.time**: A Python library to manage time delays between operations to mimic

 **human behavior.BeautifulSoup**: A Python library from bs4 for parsing HTML and XML
documents for structured data extraction (if used anywhere in extended implementations).

# EXTRACTED DATA FIELDS:

• **Job Title**:  The role or position advertised.

• **Company Name**:  The organization hiring.

• **Location**:  The city or region where the job is based.

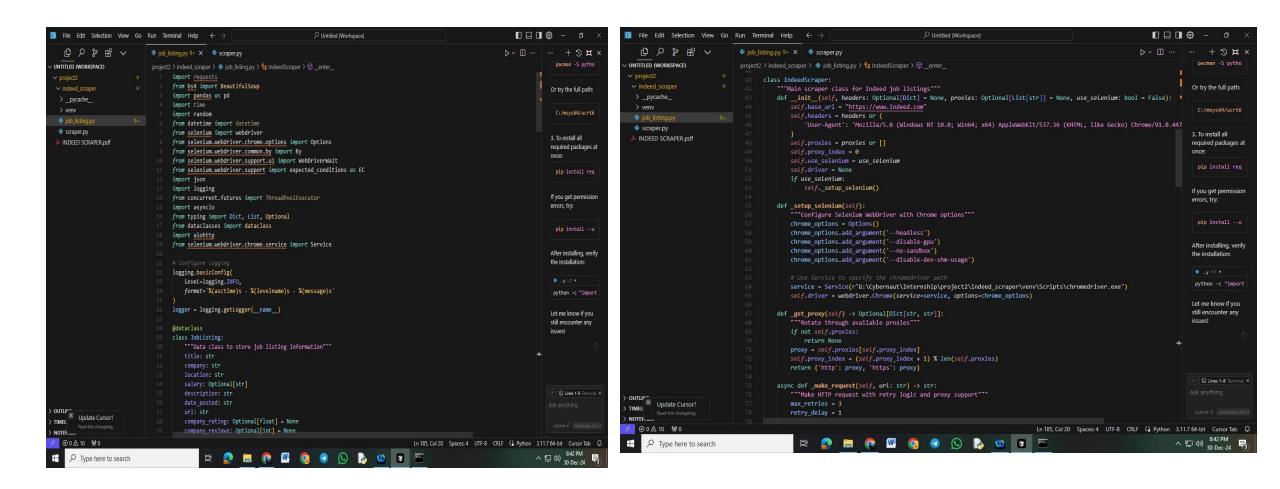• **Job Type**:  Full-time, part-time, contract, etc.

• **Salary Range**:  If available, the salary or compensation details.

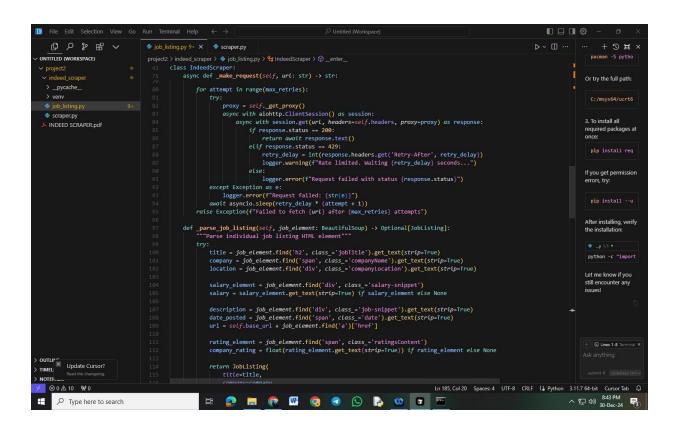• **Job Description**:  A snippet or summary of the job responsibilities.

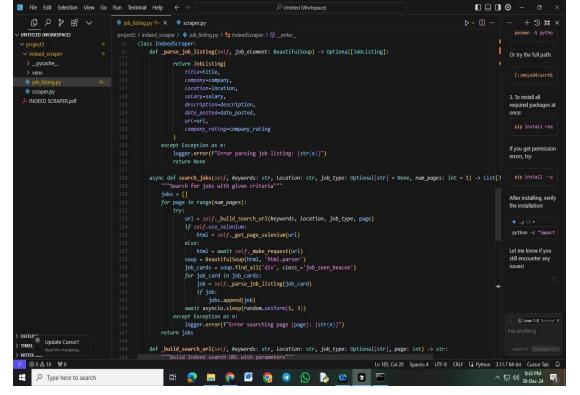• **Experience Level**:  Entry-level, mid-level, senior, etc.

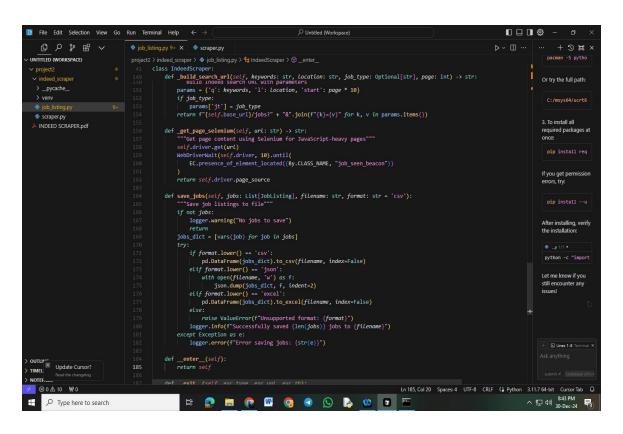• **Job Link**:  URL to the specific job listing.
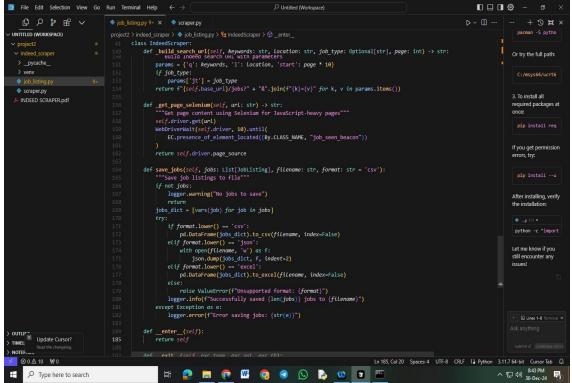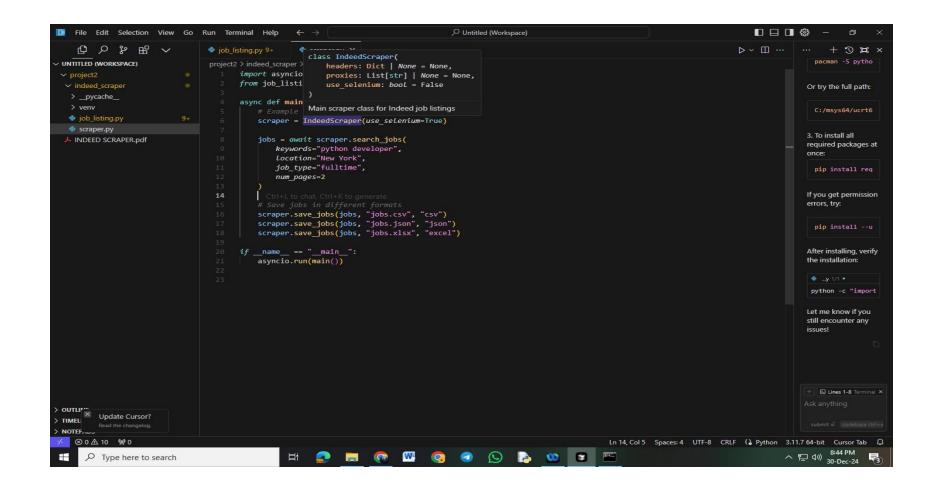
# CODE

# CODE

# CODE

# OUTPUT:

- Job Title: Python Developer | Company: Tech Innovators Inc. | Location: New York, NY | Posted: 2 days ago |

- Salary: $90,000 - $110,000 per year | Job Link: https://example.com/job/python-developer-12345

- Job Title: Backend Python Engineer | Company: Data Solutions Ltd. | Location: Remote | Posted: 5 days ago | Salary: $80,000 - $95,000 per year |

- Job Link: https://example.com/job/backend-python-engineer-67890Job Title: Junior Python Developer | Company: Future Tech Co. | Location: New York, NY | Posted: 1 day ago | Salary: $70,000 - $85,000 per year |

- Job Link: https://example.com/job/junior-python-developer-11223Job Title: Python Data Scientist | Company: Analytics World |

- Location: Brooklyn, NY | Posted: 3 hours ago | Salary: $100,000 - $130,000 per year | Job Link: https://example.com/job/python-data-scientist-44556Job Title: Lead Python Developer | Company: AI Ventures |

- Location: New York, NY | Posted: 1 week ago | Salary: $120,000 - $140,000 per year | Job Link: https://example.com/job/lead-python-developer-77889Job Title: Python Developer Intern | Company: Coding Startups |

- Location: Manhattan, NY | Posted: 2 weeks ago | Salary: Unspecified | Job Link: https://example.com/job/python-developer-intern-99001Job Title: Python Automation Engineer | Company: Smart Tech Solutions |

- Location: Remote | Posted: 4 days ago | Salary: $85,000 - $100,000 per year | Job Link: https://example.com/job/python-automation-engineer-33445

THANK YOU