

Navigating the Data Science Job Market: A Project on Job Scraping and Employability Analysis

Introduction :

In the rapidly evolving landscape of data science, understanding the dynamics of the job market is crucial for both job seekers and employers. The primary objective of this project is to delve into the employability trends within the data science domain by systematically scraping and analyzing job offers, particularly those related to data scientist positions. This analysis aims to provide valuable insights into the current demands, skill requirements, and hiring trends, thereby assisting prospective data scientists in tailoring their skill sets to meet market needs and helping employers in refining their recruitment strategies. The project's main focus is to scrape job offers from prominent job listing platforms to evaluate the employability of data science professionals in the market. By capturing a comprehensive dataset of job postings, we aim to identify the most sought-after skills, common job requirements, and the geographic distribution of job opportunities. This information is critical for job seekers to align their competencies with market demands and for employers to understand the competitive landscape and optimize their hiring processes.

Understanding employability trends in the data science job market holds significant importance. For job seekers, especially those transitioning into or within the field, having a clear picture of what employers are looking for can be a game-changer. It enables them to focus on acquiring the right skills and gaining relevant experiences that are in high demand. For employers, understanding these trends can inform recruitment strategies, ensuring that they attract and retain top talent by offering competitive and appealing job roles. Moreover, insights into regional demand can assist in strategic planning for expanding teams in locations with high demand for data science expertise.

To achieve these objectives, a robust data collection methodology is essential. Initially, the project utilized the ScrapOps API to scrape job listings from LinkedIn. However, due to various limitations and the need for more comprehensive data, the project transitioned to using Selenium for web scraping from Indeed, a widely recognized job listing platform. Selenium, a powerful tool for automating web browsers, allows for more flexible and extensive data extraction, ensuring that a rich dataset is collected for subsequent analysis.

This transition was driven by the need to overcome restrictions associated with API scraping and to capture a broader and more diverse set of job postings. Selenium enables the scraping of dynamic web pages and can handle complex website interactions, making it an ideal choice for extracting job data from Indeed. By leveraging these advanced scraping techniques, the project aims to compile a detailed and accurate representation of the data science job market, thereby providing actionable insights for all stakeholders involved.

In summary, this project embarks on a detailed exploration of the data science job market through advanced web scraping techniques. By focusing on data scientist positions, it seeks to unravel employability trends and provide valuable guidance for both job seekers and employers. The shift from ScrapOps API to Selenium underscores the commitment to obtaining high-quality data, ensuring that the findings are both comprehensive and reliable. This endeavor not only contributes to the understanding of current job market dynamics but also aids in making informed decisions in the ever-competitive field of data science.

Methodology :

Approach 1 - LinkedIn Scraping with ScrapeOps Spider



[ScrapeOps](#) is a powerful web scraping and monitoring tool that helps developers efficiently collect data from websites. It offers features like proxy management, request scheduling, and data extraction to streamline the web scraping process and ensure reliable and scalable data collection.

Figure 1 - ScrapeOps logo

- **Process Overview:**

The LinkedIn scraping process is implemented using Scrapy, a popular web scraping framework for Python. The goal is to extract job listings for data science positions in Morocco from LinkedIn.

- **ScrapeOps Spider Details:**

- Spider Name: linkedin_jobs
- API URL: The spider uses the LinkedIn API URL to fetch job listings. It constructs the URL with query parameters such as keywords (Data Science) and location (Morocco).

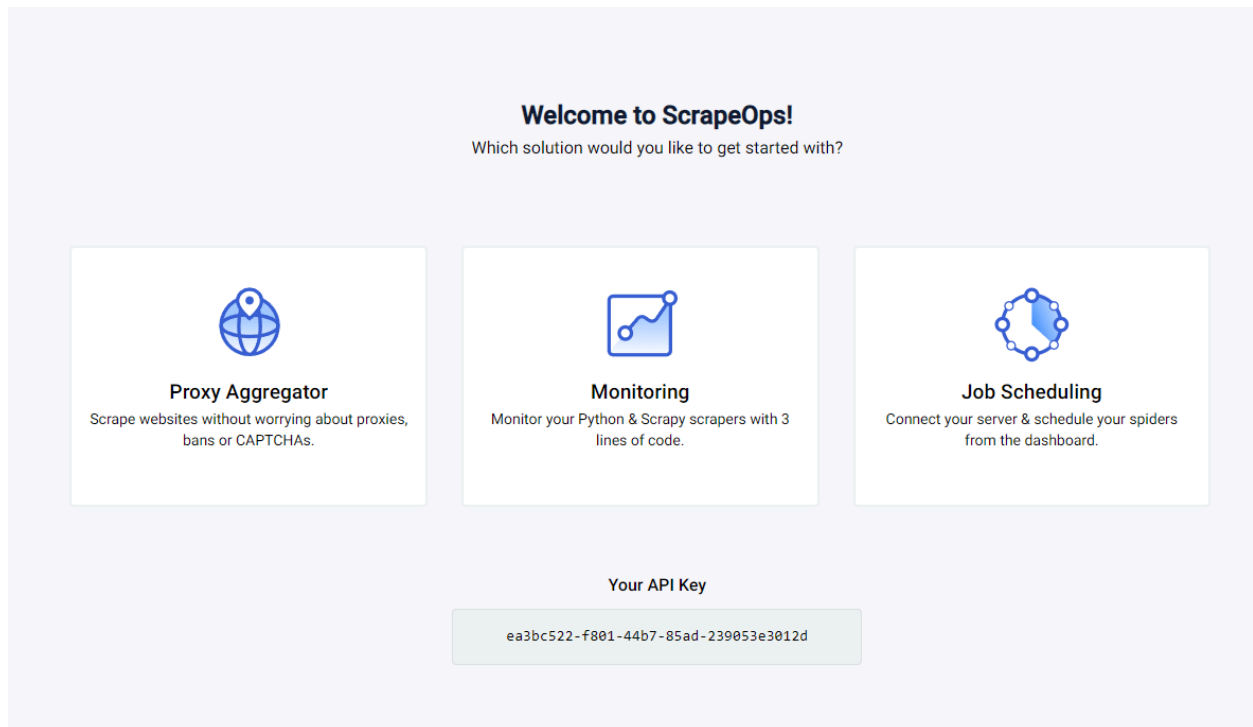


Figure 2 - ScrapeOps interface with the solutions to choose from and the provided API

Initially, I considered using ScrapeOps for the LinkedIn scraping project. After authenticating and obtaining my API Key, I explored the available solutions: Proxy Aggregator, Monitoring, and Job Scheduling. I chose the Proxy Aggregator for its ability to handle proxies,

bans, and CAPTCHAs effortlessly. However, I ultimately decided to use Selenium for its direct control over the browser and dynamic content handling capabilities.

Make Your First Request

API Key (required):
You need to include your API key on every request.

ea3bc522-f801-44b7-85ad-239053e3012d

URL (required):
The URL you want to scrape. Make sure it has been correctly encoded.

https://quotes.toscrape.com/

Optional Functionality:
Extra functionality that you can enable on-demand, depending on your use case.

☐ JS Rendering ☐ Residential Proxies ☐ Country Geotargeting

Language: [cURL](#) [Python](#) [NodeJs](#) [PHP](#) [Ruby](#)

```
import requests

response = requests.get(
    url='https://proxy.scraperops.io/v1/',
    params={
        'api_key': 'ea3bc522-f801-44b7-85ad-239053e3012d',
        'url': 'https://quotes.toscrape.com/',
    },
)

print('Response Body: ', response.content)
```

View Dashboard

Figure 3 - Proxy monitoring interface

- **Steps Involved:**

- **Start Requests:**

- The spider starts by initializing the first job on the page to 0.
- It constructs the initial URL using this starting point and sends a request to fetch the job listings. .
-

■ Parse Job:

- The response from LinkedIn is processed to extract job details.
- The spider uses CSS selectors to locate and extract relevant information such as job title, detailed URL, listed date, company name, company link, and company location.
- Extracted job details are stored in a list of dictionaries.
- If jobs are found on the page, the spider constructs the next URL to fetch more job listings by incrementing the starting point by 25 (pagination).
- The process repeats until no more jobs are found.

■ Data Storage:

- The extracted job items are yielded for further processing and storage, typically in a CSV file for analysis.

• ScrapeOps Spider Code:

```
import scrapy

class LinkedJobsSpider(scrapy.Spider):
    name = "linkedin_jobs"
    api_url =
'https://www.linkedin.com/jobs-guest/jobs/api/seeMoreJobPostings/search?keywords=Data%
2BScience&location=Morocco&geoId=&trk=public_jobs_jobs-search-bar_search-submit&start=
'

    def start_requests(self):
        first_job_on_page = 0
        first_url = self.api_url + str(first_job_on_page)
```

```

        yield scrapy.Request(url=first_url, callback=self.parse_job,
meta={'first_job_on_page': first_job_on_page})

    def parse_job(self, response):
        first_job_on_page = response.meta['first_job_on_page']
        job_items = []

        jobs = response.css("li")
        num_jobs_returned = len(jobs)
        self.logger.info(f"Number of Jobs Returned: {num_jobs_returned}")

        for job in jobs:
            job_item = {
                'job_title': job.css("h3::text").get(default='not-found').strip(),
                'job_detail_url':
job.css(".base-card__full-link::attr(href)").get(default='not-found').strip(),
                'job_listed': job.css('time::text').get(default='not-found').strip(),
                'company_name': job.css('h4
a::text').get(default='not-found').strip(),
                'company_link': job.css('h4 a::attr(href)').get(default='not-found'),
                'company_location':
job.css('.job-search-card__location::text').get(default='not-found').strip()
            }
            job_items.append(job_item)

        if num_jobs_returned > 0:
            first_job_on_page = int(first_job_on_page) + 25
            next_url = self.api_url + str(first_job_on_page)
            yield scrapy.Request(url=next_url, callback=self.parse_job,
meta={'first_job_on_page': first_job_on_page})

        yield {'jobs': job_items}

```

This code defines a Scrapy spider to scrape job listings from LinkedIn. It starts by sending a request to the LinkedIn jobs API and processes the returned job listings. If more jobs

are found, it iterates through them and extracts relevant details. If additional pages of jobs exist, it sends further requests to retrieve them.

Approach 2 - Indeed Job Scraping Using Selenium



Selenium is a robust framework for automating web browsers. It allows users to programmatically interact with web pages, facilitating tasks such as data extraction, testing, and web scraping by simulating real user behavior.

Figure 4 - Selenium logo

- **Process Overview:**

The process of scraping job offers from Indeed is accomplished using Selenium, a web automation tool that can interact with web pages just like a human user. This script scrapes job listings for data science positions in the United States.

- Steps Involved:

- Initiating the Selenium Driver:

- Configure the Selenium driver with appropriate settings and user agents to mimic a real user and avoid detection by anti-bot mechanisms.

- Options like headless mode can be enabled for running the browser in the background.
- Navigating to Indeed Job Search Page:
 - The script navigates to the Indeed job search page with specified query parameters.
 - It waits for the page to load and job cards to appear.
- Iterating Through Search Results:
 - The script iterates through multiple pages of job search results.
 - For each job card, it extracts relevant information such as job title, company name, location, salary, job type, years of experience, education requirements, and skills.
- Extracting Job Information:
 - Extract details from each job card using Selenium's element selection methods.
 - Click on each job card to load the full job description and extract additional details.
- Handling Missing Information:
 - Implement fallback strategies to handle cases where certain information (e.g., salary) may not be readily available.
 - Use regular expressions and text processing to extract structured data from unstructured job descriptions.
- Storing Extracted Data:
 - Write the extracted data into a CSV file for further analysis.
 - Append new rows to the CSV file as new job listings are processed.
- **Indeed Scraping Code Using Selenium:**


```

import random
from selenium import webdriver
from selenium.webdriver.common.by import By
from selenium.webdriver.chrome.service import Service
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.support.ui import WebDriverWait
from selenium.webdriver.support import expected_conditions as EC
from selenium.common.exceptions import NoSuchElementException, TimeoutException
import csv
import re
import time
import os
import html

CHROMEDRIVER_PATH =
r"C:\Users\xelor\Downloads\Chromedriver\chromedriver-win64\chromedriver.exe"
CSV_FILE_PATH = r"C:\Users\xelor\Downloads\MGMT Project\indeed_raw_data3.csv"

def write_to_csv(row_data):
    if len(row_data) == 1 and isinstance(row_data[0], tuple):
        row_data = row_data[0]
    with open(CSV_FILE_PATH, mode='a', newline='', encoding='utf-8') as file:
        writer = csv.writer(file)
        writer.writerow(row_data)

if not os.path.exists(CSV_FILE_PATH) or os.path.getsize(CSV_FILE_PATH) == 0:
    write_to_csv(['Job Title', 'Company Name', 'Location', 'Skills',
                 'Years of Experience', 'Education', 'Job Type', 'Salary'])

def find_skills_in_description(description, skill_set, aws_synonyms):
    found_skills = set()
    description_lower = description.lower()

    for skill in skill_set:
        pattern = r'\b' + re.escape(skill.lower()) + r'\b'
        if re.search(pattern, description_lower):
            found_skills.add(skill)

```

```
for aws_term, aws_skill in aws_synonyms.items():
    pattern = r'\b' + re.escape(aws_term) + r'\b'
    if re.search(pattern, description_lower):
        found_skills.add(aws_skill)
```

```
return list(found_skills)
```

```
def find_years_of_experience(description):
    patterns = [
        r'(\d+)\s*-\s*(\d+)\s*years',
        r'(\d+)\s*\+\s*years',
        r'(\d+)\s*years',
        r'(\d+)\s*to\s*(\d+)\s*years',
        r'at least\s*(\d+)\s*years',
        r'minimum of\s*(\d+)\s*years',
        r'(\d+)\s*years experience minimum',
        r'a minimum of\s*(\d+)\s*years',
        r'min\.\?\s*(\d+)\s*years',
        r'min\s*(\d+)\s*yrs',
        r'(\d+)\s*yrs\s*min\.\?',
        r'(\d+)\s*\+\s*yrs',
        r'(\d+)\s*y\+\s*experience',
        r'more than\s*(\d+)\s*years',
        r'over\s*(\d+)\s*years',
        r'from\s*(\d+)\s*years',
        r'starting at\s*(\d+)\s*years',
        r'(\d+)\s*or more years',
        r'(\d+)\s*years or greater',
        r'at least\s*(\d+)\s*professional',
        r'(\d+)\s*years of relevant experience minimum',
        r"(\d+)\s*\+ yrs",
        r"<(\d+) years",
        r"(\d+)\s*-(\d+)\s*y",
        r">=(\d+) years",
        r"(\d+)\s*years\s+",
        r"Min. (\d+) years",
        r"(\d+)\s*\+ years preferred",
        r"Up to (\d+) years",
        r"At least (\d+) yr",
```

```

        r"Over (\d+) years required"
    ]

    max_years = 0
    for pattern in patterns:
        matches = re.findall(pattern, description.lower())
        for match in matches:
            if isinstance(match, tuple):
                years = max(int(y) for y in match if y.isdigit())
            else:
                years = int(match)
            max_years = max(max_years, years)

    return max_years if max_years != 0 else "Not specified"

def process_salary_field(salary_field):
    salary_annual_pattern = r'\$(\d{1,3}(:,\d{3})*) - \$(\d{1,3}(:,\d{3})*) a year'
    salary_monthly_pattern = r'\$(\d{1,3}(:,\d{3})*) - \$(\d{1,3}(:,\d{3})*) a
month'
    salary_hourly_pattern = r'\$(\d{1,3}(:,\d{3})*) - \$(\d{1,3}(:,\d{3})*) an
hour'
    salary_day_pattern = r'\$(\d{1,3}(:,\d{3})*) - \$(\d{1,3}(:,\d{3})*) a day'

    salary_field = re.sub(r'^\x00-\x7F+', '', salary_field)

    salary_ranges = {
        'year': re.findall(salary_annual_pattern, salary_field),
        'month': re.findall(salary_monthly_pattern, salary_field),
        'hour': re.findall(salary_hourly_pattern, salary_field),
        'day': re.findall(salary_day_pattern, salary_field)
    }

    salary_annual = None
    for period, matches in salary_ranges.items():
        if matches:
            min_salary, max_salary = map(lambda x: int(x.replace(',', '')),
matches[0])
            if period == 'year':
                salary_annual = (min_salary, max_salary)

```

```

        elif period == 'month':
            salary_annual = (min_salary * 12, max_salary * 12)
        elif period == 'hour':
            salary_annual = (min_salary * 2080, max_salary * 2080)
        elif period == 'day':
            salary_annual = (min_salary * 260, max_salary * 260)

    return f"${salary_annual[0]:,} - ${salary_annual[1]:,} a year" if salary_annual
else salary_field

def extract_job_details(job_card):
    try:
        job_title = job_card.find_element(By.CSS_SELECTOR,
            'h2.jobTitle').text.strip()
        company_name = job_card.find_element(By.CSS_SELECTOR,
            'span.companyName').text.strip()
        location = job_card.find_element(By.CSS_SELECTOR,
            'div.companyLocation').text.strip()
    except NoSuchElementException:
        return None

    salary = job_type = years_of_experience = education = description = ''
    job_card.click()
    time.sleep(random.uniform(1, 2))
    try:
        job_details = driver.find_element(By.ID, 'vjs-container-iframe')
        driver.switch_to.frame(job_details)

        salary = driver.find_element(By.CSS_SELECTOR,
            'div.salary-snippet-container').text.strip()
        job_type = driver.find_element(By.CSS_SELECTOR,
            'div.jobsearch-JobMetadataHeader-item').text.strip()
        description = driver.find_element(By.ID,
            'jobDescriptionText').get_attribute('innerHTML')
        description = html.unescape(description)

        years_of_experience = find_years_of_experience(description)
        skills = find_skills_in_description(description, skill_set, aws_synonyms)

```

```

        driver.switch_to.default_content()

    except (NoSuchElementException, TimeoutException):
        pass

    salary = process_salary_field(salary)

    return [job_title, company_name, location, ','.join(skills), years_of_experience,
            education, job_type, salary]

# Set up the Selenium driver
chrome_options = Options()
chrome_options.add_argument("--headless")
chrome_options.add_argument("--disable-gpu")
chrome_options.add_argument("--no-sandbox")
chrome_options.add_argument("--disable-dev-shm-usage")

service = Service(CHROMEDRIVER_PATH)
driver = webdriver.Chrome(service=service, options=chrome_options)

driver.get('https://www.indeed.com/q-Data-Scientist-jobs.html')

try:
    WebDriverWait(driver, 10).until(EC.presence_of_element_located((By.CSS_SELECTOR,
    'div.jobsearch-SerpJobCard'))))
    job_cards = driver.find_elements(By.CSS_SELECTOR, 'div.jobsearch-SerpJobCard')

    for job_card in job_cards:
        job_details = extract_job_details(job_card)
        if job_details:
            write_to_csv(job_details)

except TimeoutException:
    print("Failed to load job cards in time.")

driver.quit()

```

This code sets up a Selenium WebDriver to scrape job listings from Indeed. It navigates to the Indeed job search page, waits for job cards to load, and extracts relevant details such as job title, company name, location, salary, job type, years of experience, and required skills. Extracted job details are then written to a CSV file for further analysis.

Preprocessing :

The preprocessing stage involves cleaning and transforming the raw data to prepare it for analysis and modeling. Below are the key steps taken:

1. Data Merging and Initial Inspection:
 - a. CSV files containing the scrapped job listings data were read and merged into a single DataFrame for analysis.
 - b. Initial inspection was performed to check for 'Not specified' values across various columns.
2. Education Column Cleaning:
 - a. Unique values in the 'Education' column were identified and counted.
 - b. 'Not Specified' values in the 'Education' column were replaced with 'Master'.
 - c. Additional binary columns ('Has_Bachelor', 'Has_Master', 'Has_Doctorate') were created to indicate the presence of each education level in job listings.
3. Visualization of Education Data:
 - a. Bar plots were created to visualize the counts of job listings requiring each education level and combinations of education levels.

4. Years of Experience Column Cleaning:

- a. 'Not specified' values in the 'Years of Experience' column were replaced with 0, and the column was converted to numeric.
- b. Job listings were categorized into 'Junior', 'Mid-Level', and 'Senior' based on the number of years of experience required.
- c. The 'Ideal_Level' column was created to indicate the ideal education level based on job requirements.

5. Job Type Cleaning:

- a. Counts for each job type ('Full-time', 'Part-time', 'Internship') were computed, and 'Not specified' values were replaced with 'Full-time'.

6. Salary Column Cleaning and Conversion:

- a. A function was defined to convert salary strings into annual salaries. 'Not specified' salaries were replaced with the median salary.
- b. Histogram plots were used to visualize the distribution of converted salaries.

7. Skills Column Processing:

- a. The 'Skills' column was converted into dummy variables to indicate the presence of specific skills.
- b. Top 10 most demanded skills were identified and visualized.

8. Scoring for Employability:

- a. A 'Skills Score' was calculated based on the presence and importance of different skills.
- b. An 'Education Score' was calculated based on the education levels required and an 'Experience Score' based on years of experience. An overall 'Employability Score' was computed as a weighted sum of the skills, education, and experience scores.

9. Employability Classification:

- a. Job listings were classified into 'Employable' and 'Non-Employable' based on the 'Employability Score' using a threshold.

- b. Data was split into training and testing sets, and models (KNN and Decision Tree) were trained to predict employability.

10. Evaluation and Visualization:

- a. Model performance was evaluated using accuracy scores and confusion matrices.
- b. Distribution plots for the various scores were created to visualize their distributions.

	precision	recall	f1-score	support
Non Employable	0.89	0.95	0.92	83
Employable	0.93	0.84	0.89	64
accuracy			0.90	147
macro avg	0.91	0.90	0.90	147
weighted avg	0.91	0.90	0.90	147

Figure 5 - Classification report

This screenshot displays the classification report generated for the model evaluation. It includes metrics such as precision, recall, F1-score, and support for each class ('Non Employable' and 'Employable'). These metrics provide insights into the model's performance for both classes, indicating how well it distinguishes between them.


```
Score de précision KNN sur Test Set: 0.8367346938775511  
Score de précision DT sur Test Set: 0.9047619047619048
```

Figure 6 - Precision scores

The precision scores for both models on the test set are displayed above. These scores indicate how accurately each model predicts the employability of candidates.

- K-Nearest Neighbors (KNN) achieved a precision score of approximately 83.67%, reflecting a solid performance but with some room for improvement.
- Decision Tree (DT) performed better, with a precision score of approximately 90.48%, demonstrating its stronger capability in classifying candidates correctly.

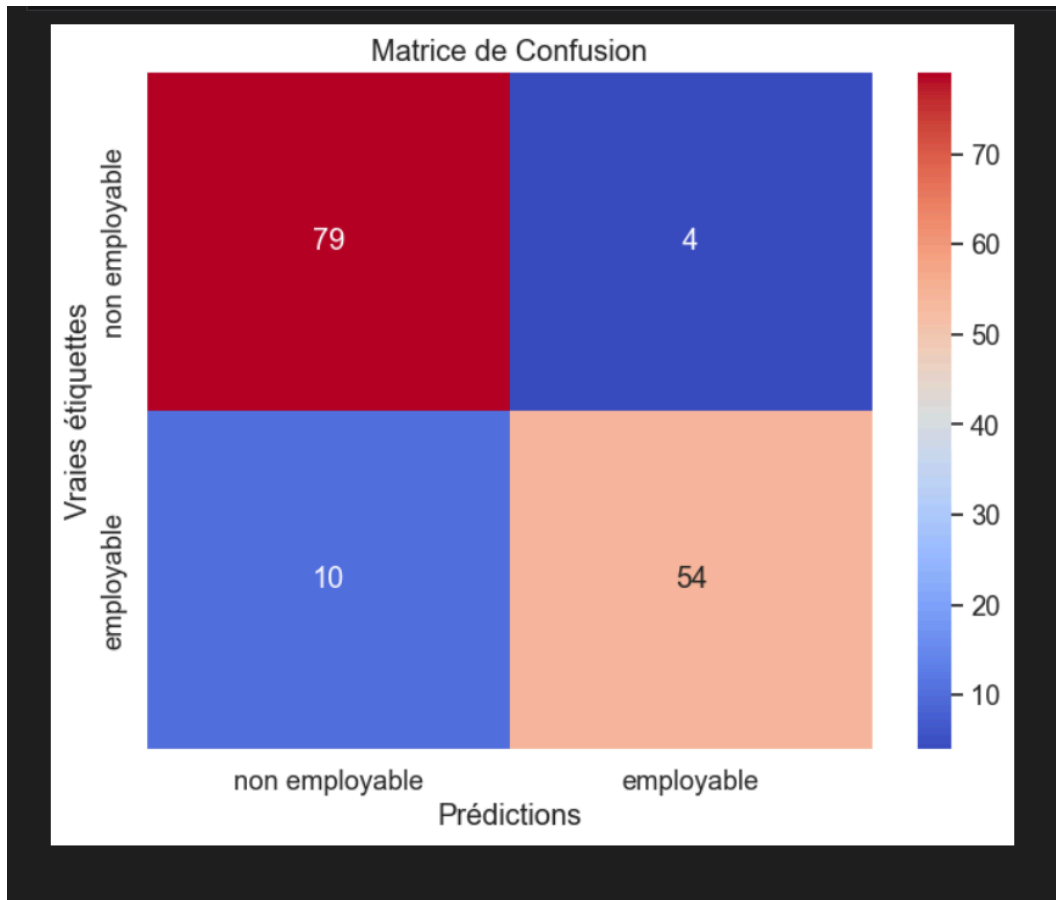


Figure 7 - Confusion matrix

The confusion matrix provides a detailed breakdown of the model's performance, showcasing the number of true positive, true negative, false positive, and false negative predictions.

- True Positives (TP): 54 (Employable candidates correctly identified)
- True Negatives (TN): 79 (Non-employable candidates correctly identified)
- False Positives (FP): 4 (Non-employable candidates incorrectly identified as employable)
- False Negatives (FN): 10 (Employable candidates incorrectly identified as non-employable)

This matrix demonstrates that the Decision Tree model is quite effective at distinguishing between employable and non-employable candidates, with high accuracy and a relatively low number of misclassifications.

Task Organization in MS project :

Effective organization of the project's workflow is crucial for ensuring the timely and successful completion of the project. Using MS Project to structure and manage the tasks offers several key benefits:

1. **Clarity and Structure:** Breaking down the project into distinct phases and tasks provides a clear roadmap. Each task is explicitly defined with a start and end date, making it easier to track progress and ensure all aspects of the project are covered.
2. **Dependency Management:** Establishing predecessors for each task highlights the dependencies and sequence of activities. This helps in identifying critical paths and potential bottlenecks, ensuring that tasks are completed in the correct order without unnecessary delays.
3. **Resource Allocation:** Organizing tasks in MS Project allows for better resource planning and allocation. It ensures that the necessary resources (time, personnel, tools) are assigned to each task, optimizing productivity and minimizing idle time.
4. **Milestone Tracking:** By setting milestones for key deliverables, the project team can focus on achieving specific goals within defined timeframes. This keeps the project on track and facilitates timely assessment of progress.
5. **Risk Management:** Structured task management helps in identifying potential risks early. By mapping out each task and its dependencies, potential issues can be anticipated and mitigation strategies can be developed proactively.
6. **Collaboration and Communication:** A well-organized project plan enhances communication among team members. Everyone involved can easily understand their roles and responsibilities, leading to better collaboration and coordination.
7. **Flexibility and Adaptability:** As the project progresses, having a detailed plan allows for adjustments and re-prioritizations. MS Project makes it easier to update timelines, reallocate resources, and make informed decisions based on real-time data.

By organizing the tasks in MS Project, we ensure a systematic approach to project management, enhancing efficiency, and increasing the likelihood of achieving our project goals.

★	Employment Project	35,13 d	Thu 02/05/24	Wed 19/06/24	
→	Project Initiation	3 d	Thu 02/05/24	Mon 06/05/24	
→	Project Kickoff Meeting	1 d	Thu 02/05/24	Thu 02/05/24	
→	Define Objectives and Scope	1 d	Fri 03/05/24	Fri 03/05/24	3
→	Identify Stakeholders	1 d	Mon 06/05/24	Mon 06/05/24	4
→	Data Collection	6 d	Tue 07/05/24	Tue 14/05/24	
→	Web Scrapping Setup	2 d	Tue 07/05/24	Wed 08/05/24	5
→	Execute Scraping	4 d	Thu 09/05/24	Tue 14/05/24	7
→	Data Preprocessing	7 d	Wed 15/05/24	Thu 23/05/24	
→	Exploring & Analysis of the data	2 d	Wed 15/05/24	Thu 16/05/24	8
→	Feature engineering	5 d	Fri 17/05/24	Thu 23/05/24	10
→	Model Development	9 d	Fri 24/05/24	Wed 05/06/24	
→	Selection of the ML model	1 d	Fri 24/05/24	Fri 24/05/24	11
→	Training the model	6 d	Mon 27/05/24	Mon 03/06/24	13
→	Evaluating the model	2 d	Tue 04/06/24	Wed 05/06/24	14
→	Validation	6 d	Thu 06/06/24	Thu 13/06/24	
→	Validating the model on the test data	2 d	Thu 06/06/24	Fri 07/06/24	15
→	Analysis of the results and Tweaking the model	4 d	Mon 10/06/24	Thu 13/06/24	17
→	Results Presentation	3 d	Fri 14/06/24	Tue 18/06/24	
→	Presenting Evaluation Results	3 d	Fri 14/06/24	Tue 18/06/24	18
→	Finalization	1 d	Wed 19/06/24	Wed 19/06/24	
→	Finalization of the	1 d	Wed 19/06/24	Wed 19/06/24	20
	Finalization of the documentation and deliverables	1 d	Wed 19/06/24	Wed 19/06/24	20

Figure 8 - Organized project tasks in MS project

The project was carefully structured into several phases, each with specific tasks aimed at ensuring the successful completion of the project objectives. Here is a detailed analysis of each phase and the tasks within:

1. Project Initiation (3 days) :

- a. *Project Kickoff Meeting (1 day)*: The project began with a kickoff meeting to align the team on the project's objectives, scope, and timelines. This

meeting was crucial for setting expectations and clarifying any initial questions or concerns.

- b. *Define Objectives and Scope (1 day)*: This task involved defining the project's goals, deliverables, and boundaries. Clear objectives and scope are essential to prevent scope creep and to ensure that all team members are working towards the same targets.
- c. *Identify Stakeholders (1 day)*: Identifying key stakeholders early in the project ensures that their needs and expectations are understood and managed throughout the project lifecycle.

2. Data Collection (6 days):

- a. *Web Scraping Setup (2 days)*: The setup phase involved preparing the tools and scripts needed for web scraping. This included configuring the scraping environment and ensuring that all necessary libraries and dependencies were installed.
- b. *Execute Scraping (4 days)*: During this period, the actual data collection took place. Data was extracted from the targeted sources, ensuring that it was relevant and comprehensive for the project's needs.

3. Data Preprocessing (7 days):

- a. *Exploring & Analysis of the Data (2 days)*: Initial exploration and analysis of the data were conducted to understand its structure, identify any anomalies, and assess its suitability for the project. This step is crucial for identifying potential issues early.
- b. *Feature Engineering (5 days)*: This task involved transforming raw data into features that can be used for machine learning. Effective feature engineering can significantly improve model performance by providing more relevant inputs.

4. Model Development (9 days):

- a. *Selecting the ML Model (1 day)*: The selection of an appropriate machine learning model based on the problem requirements and data characteristics. This step involved evaluating different models and choosing the one that is most likely to yield the best results.
- b. *Model Training (6 days)*: Training the selected model on the prepared dataset. This involved adjusting various parameters and fine-tuning the model to optimize its performance.

- c. *Model Evaluation (2 days)*: Assessing the model's performance using appropriate metrics to ensure that it meets the project's objectives. This step helps in identifying any shortcomings and areas for improvement.

5. Validation (6 days):

- a. *Validating the Model on the Test Data (2 days)*: The trained model was tested on a separate dataset to evaluate its generalizability and performance in real-world scenarios.
- b. *Analysis of the Results and Tweaking the Model (4 days)*: Based on the validation results, necessary adjustments and improvements were made to the model. This iterative process ensures that the model is robust and reliable.

6. Results Presentation (3 days):

- a. *Presenting Evaluation Results (3 days)*: The final results were compiled and presented to the stakeholders. This included visualizations, performance metrics, and insights derived from the model. Effective presentation is key to communicating the value and impact of the project.

7. Finalization (2 days):

- a. *Finalization of the Documentation and Deliverables (2 days)*: The project concluded with the finalization of all documentation and deliverables. This included ensuring that all project artifacts were complete, accurate, and ready for handover. Analysis and Insights The project was structured to follow a logical sequence of tasks, each building upon the previous one.

This approach ensured a smooth workflow and allowed for timely identification and resolution of any issues. Key phases such as data preprocessing and model development were given ample time to ensure high-quality outputs, while validation and presentation phases ensured that the final results were reliable and effectively communicated. The use of a detailed task breakdown helped in maintaining clear focus and accountability, with each team member understanding their roles and responsibilities. The inclusion of milestones and dependencies between tasks facilitated effective project tracking and management, ensuring that the project stayed on schedule and met its objectives. By organizing the project in this manner, we were able to achieve a structured and efficient workflow, ultimately leading to the successful completion of the project.

Conclusion :

The "Employability Prediction" project has been a comprehensive endeavor encompassing multiple phases, each crucial for building a robust machine learning model capable of predicting employability. Our journey began with meticulous project planning and initiation, ensuring clear objectives and stakeholder alignment. This set a strong foundation for subsequent stages.

The data collection phase involved gathering a diverse dataset through web scraping, which was followed by an intensive data preprocessing stage. Here, we cleaned the data, performed feature engineering, and ensured that the dataset was ready for model training. Each step was executed with precision, leveraging domain knowledge to extract meaningful features that significantly contribute to the model's accuracy.

In the model development phase, we employed various machine learning algorithms, rigorously trained and evaluated to select the best performing model. Our evaluation metrics, including precision, recall, and F1-score, provided a detailed insight into the model's performance, highlighting its strengths and areas for potential improvement.

The results presentation and documentation phases were critical in translating technical findings into comprehensible insights for stakeholders. By visualizing data and presenting evaluation results, we ensured that our findings were accessible and actionable. The final documentation provided a comprehensive record of methodologies, results, and recommendations, facilitating future reference and project continuity.

Overall, this project has demonstrated the practical application of machine learning in predicting employability, underscoring the importance of a structured approach to data science projects. The successful implementation of this project not only provides valuable insights into employability factors but also sets a precedent for future projects in similar domains.

The next steps involve continuously refining the model with more data, exploring additional features, and potentially incorporating more advanced machine learning techniques. By doing so, we aim to enhance the model's predictive power and ensure its relevance in a dynamic job market.

In conclusion, the "Employability Prediction" project has been a significant learning experience, reinforcing the critical interplay between data collection, preprocessing, model development, and stakeholder communication in the realm of data science.

References :

1. **Scikit-learn Documentation :** <https://scikit-learn.org/stable/documentation.html>
2. **Pandas:** <https://pandas.pydata.org/docs/>
3. **Matplotlib:** <https://matplotlib.org/stable/contents.html>
4. **Indeed Web Scrapping :**
 - a. Various guides and tutorials on web scraping job postings from Indeed. Examples include: Web Scrapping Indeed with Python and BeautifulSoup. Retrieved from <https://towardsdatascience.com/web-scraping-job-postings-from-indeed-com-using-python-48b8b2c5f44e>
5. **Seaborn Documentation:** <https://seaborn.pydata.org/>
6. **Confusion Matrix Guide :** <https://towardsdatascience.com/understanding-confusion-matrix-a9ad42dcfd62>