# Indeed Job Scraping and Analysis Report

## Task Description

Scrape all job postings from Indeed.com as of September 1, 2024, and identify all IT-related jobs.

## Initial Thoughts and Challenges Faced

- Stealth Settings: My first thought was how I could balance comprehensive data collection with efficient scraping techniques to avoid overwhelming servers or triggering rate-limiting measures.
- Javascript: While browsing Indeed before coding, I noticed they used JavaScript, so basic HTTP requests would not suffice for data collection. I used Selenium WebDriver, which mimics user interactions more effectively. While this solution offered excellent capabilities, it also introduced new challenges like slower scraping speeds and more complex error handling.
- Standardizing Title: The non-standardized labeling of IT jobs across different companies presented a unique challenge. To address this, I explored NLP techniques and implemented K-means clustering to effectively categorize and analyze the diverse job titles.
- Data quality and consistency: The job postings contained unstructured or semi-structured data, so I had to create a method to clean and standardize this. The salary information was particularly challenging due to the various formats (per hour, year, etc.). If I had more time, this would be my first improvement.

These challenges shaped my approach to the project, influencing my technical decisions and the overall design of the data pipeline. While I was able to address many of these issues, some—like the cap on the number of listings scraped—remain areas for potential future improvement.

## Design and Implementation

### Web Scraping

The core of my data collection process relied on Selenium WebDriver, configured with Chrome and enhanced with stealth settings, so I would not be detected as a bot. This approach allowed me to navigate the site to search for and extract necessary information. I developed custom functions for the search process, pagination, and data extraction.

### Data Processing and Analysis

Once collected, the raw job data was processed using Pandas for data manipulation and storage. I then employed the NLTK library and scikit-learn to identify IT-related positions. This involved tokenizing job titles, removing stop words, and applying TF-IDF vectorization to create numerical representations of the text data. From there, I used K-means clustering to group similar job titles, allowing me to identify clusters likely to contain IT-related positions.

### Cloud Integration

Thinking about scalability, I  integrated my pipeline with Google Cloud Platform where job data was stored in Google Cloud Storage. Additionally, I implemented BigQuery, allowing for efficient querying and analysis of the collected data.

### Improvements Made During Development

While developing this, I made several improvements to enhance the reliability and effectiveness of my pipeline. I added error handling and logging, ensuring that any issues during the scraping process were noted and could be debugged faster. I also added retry logic for page loading to handle timeouts, an issue I ran into while working on this project. Lastly, NLP vastly improved the number of jobs I scraped due to inconsistency in job titles.

## Limitations and Future Improvements

### Current Limitations

The current scraper only captures 15,000 job listings to prevent server overload and manage processing time. My job type extraction also gathers relatively basic information, which could be improved by adding descriptions, required qualifications, etc.

### Potential Improvements with More Time

Given more time and resources, I would explore several avenues to improve this project. Namely:

1. Implementing distributed scraping could significantly increase my data collection capacity.
2. Expanding my data analysis capabilities to include sentiment analysis of job descriptions or extraction of required skills could provide deeper insights into the job market.

On the infrastructure side, it would be helpful to set up automated scheduling for regular data updates and implement data versioning in Google Cloud Storage.