**Job Description:**
We are looking for an experienced Python Developer to build a scraper for extracting Amazon product reviews. The scraper will process URLs in the following format:

```
https://www.amazon.com/product-reviews/{dynamic-variable-asin}/?reviewerType=all_reviews&sortBy=recent&pageNumber=1
```

The scraper must securely and efficiently retrieve data, handle proxies and user agents for anonymity, and store the results in a MySQL database. For the field **"Helpful Numbers"**, the scraper should default to 0 if the value is not available. Additionally, the scraper should capture the **user link** (profile URL) for each reviewer. The scraper will run on Ubuntu 24 within a Python virtual environment.

---

## Key Responsibilities:

- **Scraper Development**:
  - Build a web scraper using Python Scrapy (or an equivalent framework) to process Amazon product review pages with dynamic ASINs.
  - Handle pagination by traversing multiple review pages until no new reviews are found.
- **Proxy and User-Agent Rotation**:
  - Integrate proxy management to assign or rotate proxies randomly with authentication (proxies will be provided).
  - Rotate user agents to emulate various devices and browsers.
- **Data Extraction**:
  - Extract the following fields from each review:
    - **Review ID**
    - **User Name**
    - **User Link** (profile URL)
    - **Review Title**
    - **Review Ratings**
    - **Created Date**
    - **Review Content**
    - **Helpful Numbers** (default to 0 if the value is unavailable)
- **Database Integration**:
  - Store extracted reviews in a MySQL database.
  - Avoid duplicate entries by checking for existing reviews based on the **Review ID**.
- **IP Protection**:
  - Prevent exposure of the server's IP address by effectively using proxies and other security measures.

## Technical Requirements:

- **Programming**: Strong proficiency in Python and frameworks like Scrapy, Selenium, or BeautifulSoup.
- **Proxies**: Expertise in proxy management with support for authentication and rotation.
- **User Agents**: Familiarity with configuring and rotating user agents.
- **Database**: Experience with MySQL, including schema design and writing efficient queries for checking duplicates.
- **Environment**: Ability to run Python scripts in a virtual environment on Ubuntu 24.
- **Condition Handling**: Implement logic to set **Helpful Numbers** to $0$ if the value is not available.
- **User Link Extraction**: Ability to extract the user's profile link for each review.

## Deliverables:

1. **Python Scraper**: A script to scrape Amazon reviews for any ASIN using the specified URL format.
2. **MySQL Integration**: A database schema to store reviews, ensuring duplicates are avoided.
3. **Documentation**: Detailed setup and usage instructions, covering the virtual environment, database configuration, and scraper execution.
4. **Security Features**: Proxy and user-agent rotation for anonymity.
5. **Data Handling**: Correct implementation of the "Helpful Numbers" field with a fallback value of $0$.
6. **User Link Extraction**: Extraction of the user's profile link for each reviewer.