

# Take-home Assignment

## Data Engineer (Crawling Focus)

**FullName: Dang Tien Dat**

**Position: Data Engineer**

### Table of Contents

1. Data Source .....	2
2. Define Data Schema: .....	2
3. Data Extraction .....	5
a. Data Extraction Techniques: .....	5
b. Data Extraction Challenges and Solutions:.....	6
4. Data Storage in MongoDB.....	6
5. Data Pipeline Orchestration .....	8
6. Instruction: .....	10
Link Source Code: <a href="https://github.com/TienDat57/qode-assignment">https://github.com/TienDat57/qode-assignment</a> .....	11

## 1. Data Source

- LinkedIn: LinkedIn is a professional networking platform widely used by professionals across various industries. It provides a wealth of information about individuals' professional backgrounds, skills, education, work experience, and endorsements.
  - o Link: <https://www.linkedin.com/>
- JobGo: JobGo is a popular job search platform that aggregates job listings from various sources. While primarily known for job postings, JobGo also contains valuable data about candidates, particularly those actively seeking employment
  - o Link: <https://employer.jobsgo.vn/>

## 2. Define Data Schema:

The data schema is implemented using Python's Scrapy framework, utilizing Scrapy items to define the structure of the data to be scraped from web sources. The PersonProfileItem class encapsulates the various attributes associated with candidate profiles. Below is a breakdown of the fields included in the schema:

- **id**: Unique identifier for the candidate profile.
- **fullname**: Full name of the candidate.
- **title**: Current job title or professional designation.
- **location**: Geographic location of the candidate.
- **industry**: Industry or sector in which the candidate operates.
- **current\_company**: Name of the candidate's current employer or company.
- **past\_company**: List of past employers or companies the candidate has worked for.

- **education:** Educational qualifications and institutions attended by the candidate.
- **profile\_link:** URL link to the candidate's profile.
- **summary:** Summary or bio section describing the candidate's background and expertise.
- **website:** Personal or professional website associated with the candidate.
- **skill:** List of skills and expertise possessed by the candidate.
- **email:** Email address of the candidate.
- **image:** URL link to the candidate's profile picture.
- **experience:** Detailed work experience, including job roles, responsibilities, and dates.
- **birthday:** Date of birth or age of the candidate.
- **phone:** Contact phone number of the candidate.

```
Tien Dat Dang, 1 hour ago | 1 author (Tien Dat Dang)
6  class PersonProfileItem(Item):
7      id = Field()
8      fullname = Field()
9      title = Field()
10     location = Field()
11     industry = Field()
12     current_company = Field()
13     past_company = Field()
14     education = Field()
15     profile_link = Field()           Tien Dat Dang,
16     summary = Field()
17     website = Field()
18     skill = Field()
19     email = Field()
20     image = Field()
21     experience = Field()
22     birthday = Field()
23     phone = Field()
```

Figure 1. Schema

#### NOTE: Additional Schema in /schema Directory:

- The schema located in the /schema directory is described as more comprehensive and specific.
- It likely includes a broader range of attributes and may offer more detailed fields for capturing nuanced information about candidates.
- The schema may be structured to accommodate specific requirements or preferences of the organization or project, providing a tailored approach to candidate profiling.

### 3. Data Extraction

#### a. Data Extraction Techniques:

##### 1. Web Scraping with Scrapy:

- Scrapy, a Python web crawling and scraping framework, was utilized for extracting data from web sources.
- Custom spiders were developed to navigate through the structure of web pages, locate relevant candidate information, and extract desired attributes.
- XPath selectors and CSS selectors were employed to target specific HTML elements containing candidate data, such as profile details, job listings, and educational backgrounds.

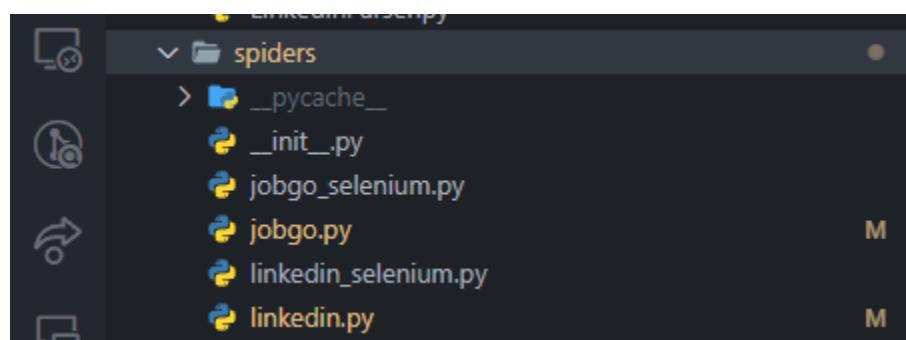
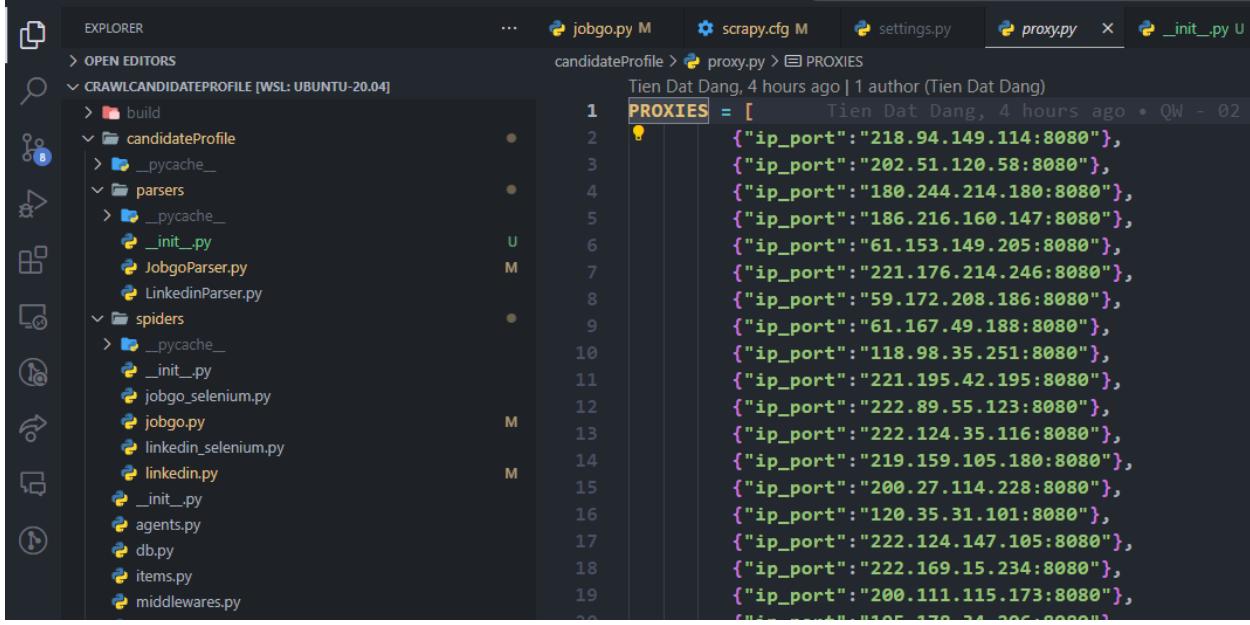


Figure 2. List spider

##### 2. Handling Anti-Scraping Measures:

- Both LinkedIn and JobGo employ anti-scraping measures to prevent automated access to their platforms.
- To bypass these measures, techniques such as user-agent rotation, IP rotation, and request throttling were implemented.
- Randomized delays between requests were introduced to simulate human-like browsing behavior and avoid triggering anti-scraping mechanisms.



The screenshot shows a code editor interface with several tabs at the top: jobgo.py, scrapy.cfg, settings.py, proxy.py (which is the active tab), and \_\_init\_\_.py. The proxy.py tab displays a Python code snippet defining a list of proxies. The code is as follows:

```

1 PROXIES = [
2     {"ip_port": "218.94.149.114:8080"}, 
3     {"ip_port": "202.51.120.58:8080"}, 
4     {"ip_port": "180.244.214.180:8080"}, 
5     {"ip_port": "186.216.160.147:8080"}, 
6     {"ip_port": "61.153.149.205:8080"}, 
7     {"ip_port": "221.176.214.246:8080"}, 
8     {"ip_port": "59.172.208.186:8080"}, 
9     {"ip_port": "61.167.49.188:8080"}, 
10    {"ip_port": "118.98.35.251:8080"}, 
11    {"ip_port": "221.195.42.195:8080"}, 
12    {"ip_port": "222.89.55.123:8080"}, 
13    {"ip_port": "222.124.35.116:8080"}, 
14    {"ip_port": "219.159.105.180:8080"}, 
15    {"ip_port": "200.27.114.228:8080"}, 
16    {"ip_port": "120.35.31.101:8080"}, 
17    {"ip_port": "222.124.147.105:8080"}, 
18    {"ip_port": "222.169.15.234:8080"}, 
19    {"ip_port": "200.111.115.173:8080"}, 
20    {"ip_port": "105.178.34.206:8080"}]

```

Figure 3. List Proxy

## b. Data Extraction Challenges and Solutions:

### CAPTCHA Challenges:

- LinkedIn and JobGo occasionally present CAPTCHA challenges to verify the legitimacy of user access, hindering automated scraping.

### Solutions:

- CAPTCHA solving services or manual intervention were employed as temporary solutions to bypass CAPTCHA challenges.
- Utilizing proxy services to rotate IP addresses and avoid triggering CAPTCHA checks.

## 4. Data Storage in MongoDB

The candidate data stored in MongoDB is organized in collections, each representing a distinct entity or type of information. The schema design aims to maintain flexibility while accommodating the diverse attributes associated with candidate profiles. Here is an overview of the database schema:

```
Tien Dat Dang, 2 hours ago | 1 author (Tien Dat Dang)
from candidateProfile import settings
import pymongo

Tien Dat Dang, 2 hours ago | 1 author (Tien Dat Dang)
class MongoDBClient(object):
    def __init__(self, col, index=None):
        connection = pymongo.MongoClient(settings.MONGODB_URI)
        self.name_col = col
        self.db = connection[settings.MONGODB_DB]
        self.collection = self.db[col]
        if index:
            self.collection.create_index(index, unique=True)

    def refresh_collection(self):
        self.collection.drop()
        self.collection = self.db[self.name_col]

    def get_collection(self):
        return self.collection

    def _walk(self):
        skip = 0
        limit = 1000
        hasMore = True
        while hasMore:
            res = self.collection.find(skip=skip, limit=limit)
            hasMore = (res.count(with_limit_and_skip=True) == limit)
            for x in res:
                yield x
            skip += limit

    def walk(self):
        docs = []
        for doc in self._walk():
            docs.append(doc)
        return docs
```

Figure 4. Declare connection with Mongo

## 5. Data Pipeline Orchestration

Zyte provides powerful tools for web scraping, data extraction, and pipeline orchestration, enabling the seamless integration of data from multiple sources into a unified workflow. The report discusses the architecture of the data pipeline, the role of Zyte in orchestrating the pipeline, and the benefits of using Zyte for managing and automating data extraction tasks.

Link: <https://www.zyte.com/>

Zyte's scheduling and task management features enable the automation and coordination of data extraction tasks.

### Add Periodic Job

**Spiders** **Choose Month**

Jobgo  Every month 

**Priority** **Choose Day Of Week**

Highest  Every day 

**Description** **Choose Day Of Month**

Add Periodic Job:  Every day 

**Tags** **Choose Hour**

run-jobgo  01:00 

Type the tag names in the field above, use **comma** (",") or **space (" ")** or press **enter** to define them.

**Arguments** 

59 

**Cancel** **Save**

Figure 5. Add schedule Job JobGo

### Add Periodic Job

<p><b>Spiders</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">linkedin <span style="color: red;">X</span></div>	<p><b>Choose Month</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">Every month <span style="float: right;">▼</span></div>
<p><b>Priority</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">High <span style="float: right;">▼</span></div>	<p><b>Choose Day Of Week</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">Every day <span style="float: right;">▼</span></div>
<p><b>Description</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">Add Periodic Job: <u>linkedin</u></div>	<p><b>Choose Day Of Month</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">Every day <span style="float: right;">▼</span></div>
<p><b>Tags</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">run-linkedin <span style="color: red;">X</span></div> <p>Type the tag names in the field above, use <b>comma</b> (",") , <b>space (" ")</b> or press <b>enter</b> to define them.</p>	<p><b>Choose Hour</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">01:00 <span style="float: right;">▼</span></div>
<p><b>Arguments</b> <span style="color: green;">+</span></p>	<p><b>Choose Minute</b></p> <div style="border: 1px solid #ccc; padding: 5px; border-radius: 10px; width: 90%;">55 <span style="float: right;">▼</span></div>
<span style="border: 1px solid #ccc; border-radius: 10px; padding: 5px 10px; color: #000; margin-right: 10px;">Cancel</span> <span style="background-color: #e64a89; color: white; border: 1px solid #e64a89; border-radius: 10px; padding: 5px 10px; font-weight: bold;">Save</span>	

Figure 6. Add schedule Job Linkedin

Periodic Jobs / Spiders																																													
Periodic Jobs																																													
JOBS																																													
<a href="#">Dashboard</a>																																													
<a href="#">Periodic Jobs</a>																																													
<a href="#">SPIDERS</a>																																													
<a href="#">Dashboard</a>																																													
<a href="#">Code &amp; Deploy</a>																																													
<a href="#">Settings</a>																																													
<a href="#">PROJECT</a>																																													
<a href="#">Usage Stats</a>																																													
<a href="#">Activity</a>																																													
<a href="#">Members</a>																																													
<a href="#">Settings</a>																																													
<a href="#">ADDONS</a>																																													
<a href="#">Addons Setup</a>																																													
<a href="#">Add periodic job</a>																																													
<a href="#">Watch</a> ▾																																													
<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th style="width: 20px;"></th> <th style="width: 10px;">Month</th> <th style="width: 10px;">Day of Month</th> <th style="width: 10px;">Day of Week</th> <th style="width: 10px;">Hour</th> <th style="width: 10px;">Spiders</th> <th style="width: 10px;">Priority</th> <th style="width: 10px;">Arguments</th> <th style="width: 10px;">Tags</th> <th style="width: 10px;">Description</th> <th style="width: 10px;">Enabled</th> <th style="width: 10px;">Actions</th> </tr> </thead> <tbody> <tr> <td><input type="checkbox"/></td><td>Every month</td><td>Every day</td><td>Every day</td><td>0159</td><td>Jobgo</td><td>Highest</td><td></td><td>run-jobgo</td><td>Add Periodic Job</td><td><input checked="" type="checkbox"/></td><td><span style="color: green;">Edit</span></td></tr> <tr> <td><input type="checkbox"/></td><td>Every month</td><td>Every day</td><td>Every day</td><td>0155</td><td>linkedin</td><td>High</td><td></td><td>run-linkedin</td><td>Add Periodic Job</td><td><input checked="" type="checkbox"/></td><td><span style="color: green;">Edit</span></td></tr> </tbody> </table>											Month	Day of Month	Day of Week	Hour	Spiders	Priority	Arguments	Tags	Description	Enabled	Actions	<input type="checkbox"/>	Every month	Every day	Every day	0159	Jobgo	Highest		run-jobgo	Add Periodic Job	<input checked="" type="checkbox"/>	<span style="color: green;">Edit</span>	<input type="checkbox"/>	Every month	Every day	Every day	0155	linkedin	High		run-linkedin	Add Periodic Job	<input checked="" type="checkbox"/>	<span style="color: green;">Edit</span>
	Month	Day of Month	Day of Week	Hour	Spiders	Priority	Arguments	Tags	Description	Enabled	Actions																																		
<input type="checkbox"/>	Every month	Every day	Every day	0159	Jobgo	Highest		run-jobgo	Add Periodic Job	<input checked="" type="checkbox"/>	<span style="color: green;">Edit</span>																																		
<input type="checkbox"/>	Every month	Every day	Every day	0155	linkedin	High		run-linkedin	Add Periodic Job	<input checked="" type="checkbox"/>	<span style="color: green;">Edit</span>																																		
<span style="border: 1px solid #ccc; border-radius: 10px; padding: 5px 10px; color: #000;">Delete</span>																																													

Figure 7. List Job created.

Figure 8. Job Dashboard

Figure 8, this showed job next and job running and number job is complete.

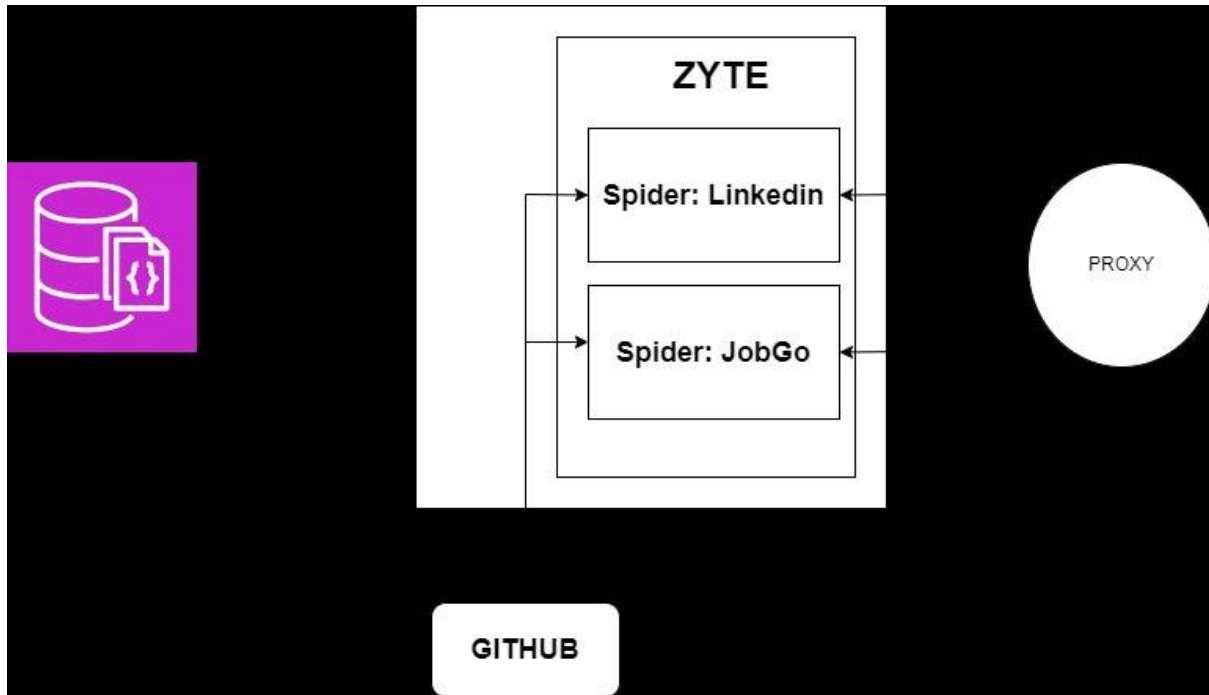


Figure 9. Architecture data pipeline

## 6. Instruction:

At file setttings.py, fill information necessary to here:

```

1 import logging
2 import datetime
3
4 logging.basicConfig(format='%(asctime)s - %(name)s - %(levelname)s - %(message)s', level=logging.DEBUG)
5 logger = logging.getLogger(__name__)
6
7 LOG_FILE = 'logs/LOG_' + datetime.datetime.now().strftime("%Y-%m-%d") + '.log'
8 LOG_LEVEL = 'DEBUG'
9
10 BOT_NAME = "candidateProfile"
11
12 SPIDER_MODULES = ["candidateProfile.spiders"]
13 NEWSPIDER_MODULE = "candidateProfile.spiders"
14
15 DOWNLOADER_MIDDLEWARES = {
16     'candidateProfile.middlewares.CustomHttpProxyMiddleware': 543, FILL HERE
17     'candidateProfile.middlewares.CustomUserAgentMiddleware': 545,
18 }
19
20 # NOTE: Fill account and password for LinkedIn and JobGo
21 LINKEDIN_ACCOUNT = 'dangtiendat1135@gmail.com' ←
22 PASSWORD_CRAWL =
23
24 JOBGO_ACCOUNT = 'dangtiendat1135@gmail.com'
25
26 # ITEM_PIPELINES = {
27 #     "candidateProfile.pipelines.MongoDBPipeline",
28 # }
29
30 # NOTE: Fill in the following fields with your MongoDB connection details
31 MONGODB_URI = 'mongodb+srv://TienDat57:<password>@qode.addyuj.mongodb.net/' ←
32 MONGODB_PORT = 27017
33 MONGODB_DB = 'qodeworld'
34 MONGODB_COLLECTION = 'candidateProfile'
35 MONGODB_UNIQ_KEY = '_id'

```

## Run local:

Requirement:

- Installed Anaconda
- Python

Open shell:

- git clone <https://github.com/TienDat57/qode-assignment>

Install environment:

- pip install -r requirements.txt

Run spider:

- scrapy crawl linkedin
- scrapy crawl Jobgo

**Link Source Code: <https://github.com/TienDat57/qode-assignment>**