

# Scraper Service-Linked In

## Scraper.py Functionality :

- db\_connect(properties):

This function takes properties which basically contains database parameters loaded from parameters.json file.

```
properties = open('parameters.json')
data = json.load(properties)
server_name = data['server_name']
user_name = data['user_name']
password = data['password']
db_name = data['db_name']
connection = mysql.connector.connect(host=server_name,
                                     database=db_name,
                                     user=user_name,
                                     password=password)

return connection
```

Returns: database connection object

- get\_all\_skills(connection)

```
def get_all_skills(connection):
    sql_select_Query = "select DISTINCT skill_id,skill_title from skill_master"
    cursor=connection.cursor()
    cursor.execute(sql_select_Query)
    records=cursor.fetchall()
    all_skills={}
    for row in records:
        all_skills[row[0]]=row[1]
    print("All skills",all_skills)
    return all_skills
```

- get\_resume\_skills(connection)

```
def get_resume_skills(connection):
    sql_select_Query2="select resume_id,skill_id from resume_skills where is_active=1"
    cursor=connection.cursor()
    cursor.execute(sql_select_Query2)
    records2=cursor.fetchall()
    resume_skills={}
    for row in records2:
        if(row[0]) in resume_skills:
            resume_skills[row[0]].append(row[1])
        else:
            resume_skills[row[0]]=row[1]
    return resume_skills
```

- get\_emailing\_list(connection):

```
def get_emailing_list(connection):
    email_id_list={}
    sql_select_Query3="SELECT resume_id,user_email from user_master um join user_resume ur on um.user_id=ur.user_id"
    cursor=connection.cursor()
    cursor.execute(sql_select_Query3)
    records_email=cursor.fetchall()
    for row in records_email:
        email_id_list[row[0]]=row[1]
    print("Resume id and email id",email_id_list)
    return email_id_list
```

Argument: connection object returned by db\_connect function()

Returns: Returns a dictionary containing a particular resume\_id and the corresponding email id of the user to whom the resume belongs to.

- get\_job\_description(keyword,no\_of\_jobs\_to\_retrieve,data):

Arguments:

1. Keyword: takes the search query i.e which job profile is the user looking for to be used as filter later.
2. No\_of\_jobs\_to\_retrieve: Takes the number of jobs user wants to search and retrieve from each site .
3. data basically has information from parameters.json which has the password for the dummy linked in account.

Execution:

1. Login in to a dummy linked in profile.

Code snippet

```
browser.get('https://www.linkedin.com/checkpoint/rm/sign-in-another-account?fromSignIn=true&trk=guest_homepage-basic_nav-header-signin')
username_ip=browser.find_element_by_id('username')
username_ip.send_keys(username)
pwd_ip=browser.find_element_by_id('password')
pwd_ip.send_keys(pwd)
sign_in_button=browser.find_element_by_xpath("//button[@data-litms-control-urn='login-submit']")
sign_in_button.click();
```

2. After logging in the web browser takes us to the home page from which we are redirecting the browser via selenium to go the job listing page.

```
##### traverse to job listing page #####
browser.get('https://www.linkedin.com/jobs/jobs-in-raleigh-nc?trk=homepage-basic_intent-module-jobs&position=1&pageNum=0')
```

3. At the job listing page we enter the search query which is stored in our keyword variable.

```
weblement = WebDriverWait(browser, 10000).until(
    EC.presence_of_element_located((By.XPATH, "//input[contains(@id,'jobs-search-box-keyword-id')]"))
)
job_description=browser.find_element_by_xpath("//input[contains(@id,'jobs-search-box-keyword-id')]").send_keys(searchquery)
#inserting job filter value
search_button=browser.find_element_by_class_name("jobs-search-box__submit-button")
search_button.click()
time.sleep(3)#give time to load search query results
```

4. After the search query is executed we scroll done to the end of the job listings page so that all possible job listings on the current page can be seen. Currently we are only retrieving jobs from the first page.

```
recentList = browser.find_elements_by_xpath("//section[@aria-label='pagination']")
for list in recentList :
    browser.execute_script("arguments[0].scrollIntoView();", list )
time.sleep(5)
```

5. Now we retrieve all the job links from the job listing panel and store them in a list

```
job_cards=browser.find_elements_by_xpath("//a[@class='disabled ember-view job-card-container__link job-card-list__title']")
href_arr=[]
for i in job_cards:
    href_arr.append(i.get_attribute("href"))
#print(Len(href_arr))
```

6. Now for each of the job listing from that list we scrape all the relevant data such as responsibilities, skills and so on (usually written in list tags on each job posting) and storing it in a list. After that, we are converting all the entries in that list to a single string variable

```
final={}
listele=[]
for url in href_arr:
    browser.get(url)
    time.sleep(5)
    show_more_button=browser.find_element_by_xpath("//button[contains(@aria-label,'Click to see more description')]")
    show_more_button.click()
    list_ele=browser.find_elements_by_xpath("//article//li")
    #####for each job listing loop through all list items and add the text#####
    data=[]
    datastr=""
    for li in list_ele:
        data.append(li.text)
    for val in data:
        datastr+=val + " "
    time.sleep(5)
    count+=1
    if(count==no_of_jobs_to_retrieve):
        break
    final[url]=datastr
```

7. This is then passed to the get\_user\_id\_to\_list\_of\_job\_ids() function.

Returns:

Final list containing email\_id and their corresponding job links is returned.

Email service

We are using the smtp protocol to send email containing the job links to each of the email id

```

1 port = 587
2 smtp_server = "smtp.gmail.com"
3 login = "srijas.alerts@gmail.com"
4 password = ""
5 sender = "srijas.alerts@gmail.com"
6 for key in final_result:
7     if key in email_id_list:
8         receiver = email_id_list[key]
9         print(receiver)
10        msg = MIMEText()
11        msg['From'] = sender
12        msg['To'] = receiver
13        msg['Subject'] = 'Job Lisitngs'
14        body = ""Hi, \n PFA the attached list of jobs that match your resume \n""
15        temp_str = ""
16        list_of_curr_links = final_result[key]
17        counter = 1
18        for link in list_of_curr_links:
19            temp_str += (str(counter) + ". " + link + '\n')
20            counter += 1
21        body += temp_str
22        msg.attach(MIMEText(body, 'plain'))
23        text = msg.as_string()
24
25        try:
26            server = smtplib.SMTP(smtp_server, port)
27            server.connect(smtp_server,port)
28            server.ehlo()
29            server.starttls()
30            server.ehlo()
31            server.login(login, password)
32            server.sendmail(sender, receiver, text)
33            server.quit()
34            print('Sent')
35        except (gaierror, ConnectionRefusedError):
36            print('Failed to connect to the server. Bad connection settings?')
37        except smtplib.SMTPServerDisconnected as e:
38            print('Failed to connect to the server. Wrong user/password?')
39            print(str(e))
40        except smtplib.SMTPException as e:
41            print('SMTP error occurred: ' + str(e))

```

## keyword\_extraction\_modules.py

This file uses four different functions:

`get_list_of_matched_skills(description, total_skills)`

This function takes the description and matches with skills from the skills table. It returns the list of matches skills

```

def get_list_of_matched_skills(description, total_skills):
    list_of_skills_matched = []
    description = description.upper()
    desc_list = description.split(" ")
    for skill in total_skills:
        if (total_skills[skill].upper() in desc_list) or ((total_skills[skill].upper() + ".") in desc_list) or (total_skills[skill].upper() in desc_list):
            list_of_skills_matched.append(skill)
    #print(list_of_skills_matched)
    return list_of_skills_matched

```

get\_dict\_with\_list\_of\_skills\_from\_description(links\_description\_dict, total\_skills)

This function creates a dictionary of job description and skills using the get\_list\_of\_matched\_skills function.

```
def get_dict_with_list_of_skills_from_description(links_description_dict, total_skills):
    job_desc_link_and_skills_dict = dict()
    for link in links_description_dict:
        description = links_description_dict[link]
        skills_matched_in_curr_description = get_list_of_matched_skills(description, total_skills)
        job_desc_link_and_skills_dict[link] = skills_matched_in_curr_description
    return job_desc_link_and_skills_dict
```

match\_both\_lists(list\_of\_skills\_in\_resume, list\_of\_skills\_in\_description, threshold, total\_skill\_count)

This function returns true if the % of skills in the resume matches skills in the job description is greater than the threshold.

```
def match_both_lists(list_of_skills_in_resume, list_of_skills_in_description, threshold, total_skill_count):
    count = 0
    for skill in list_of_skills_in_description:
        if skill in list_of_skills_in_resume:
            count += 1
    if ((count / total_skill_count) * 100 >= threshold): return True
    return False
```

get\_user\_id\_to\_list\_of\_job\_ids(resume\_skills\_dict, links\_description\_dict, db\_connection, total\_skills, threshold)

This function returns the user ids and the corresponding job links that matches with the skills in the resume.

```
def get_user_id_to_list_of_job_ids(resume_skills_dict, links_description_dict, db_connection, total_skills, threshold):
    result_dict = dict()
    job_desc_link_and_skills_dict = get_dict_with_list_of_skills_from_description(links_description_dict, total_skills)
    for curr_resume in resume_skills_dict:
        list_of_skills_in_resume = resume_skills_dict[curr_resume]
        for link in job_desc_link_and_skills_dict:
            list_of_skills_in_description = job_desc_link_and_skills_dict[link]
            if match_both_lists(list_of_skills_in_resume, list_of_skills_in_description, threshold, len(total_skills)):
                if curr_resume in result_dict:
                    curr_list = result_dict[curr_resume]
                else:
                    curr_list = []
                curr_list.append(link)
            result_dict[curr_resume] = curr_list
    return result_dict
```

