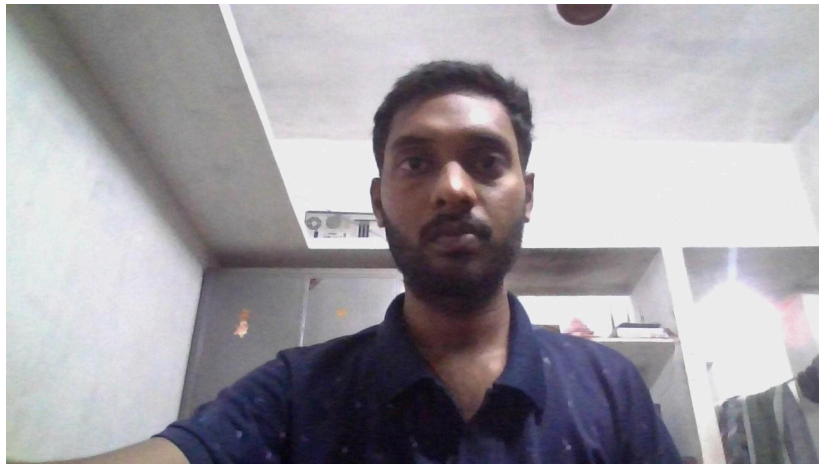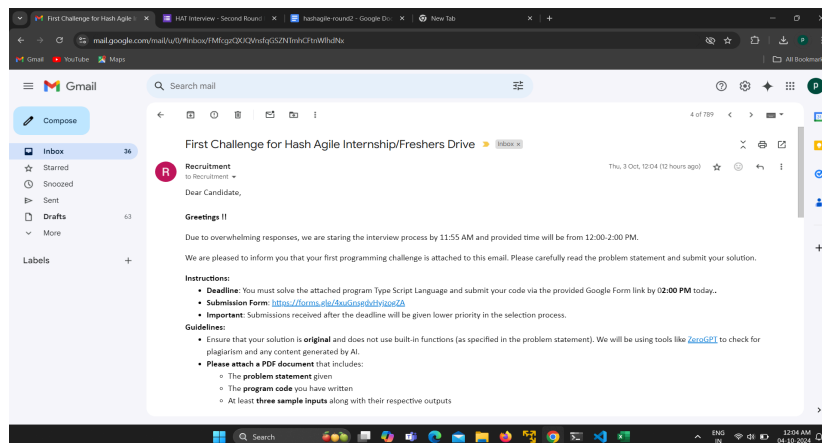# Hashagile Round-2 [submission]

A. Full Name: Venkata Ramprasad Pade
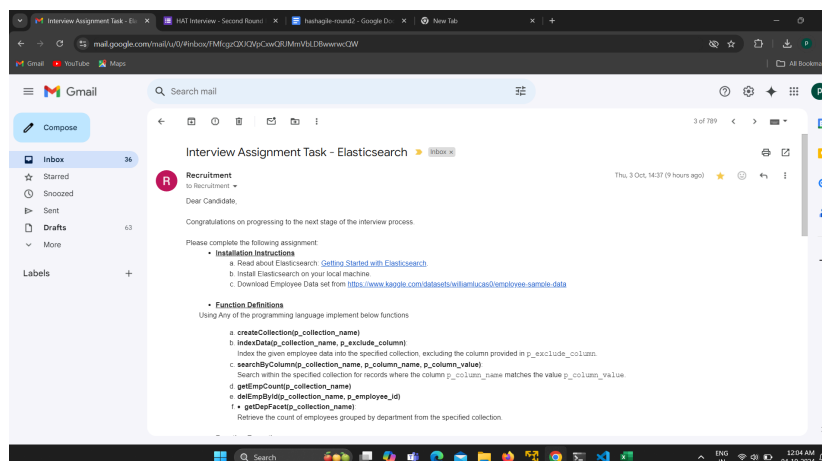
B. Selfie:



C. First Round Mail:



D. Second Round Mail



E. Github Url for first round:
https://github.com/ramprasad-13/hashagile-round1

F. Github Url for second round:
https://github.com/ramprasad-13/hashagile-round2

# Problem Statement

- **Installation Instructions**
    1. Read about Elasticsearch: [Getting Started with Elasticsearch](#).
    2. Install Elasticsearch on your local machine.
    3. Download Employee Data set from
       [https://www.kaggle.com/datasets/williamlucas0/employee-sample-data](https://www.kaggle.com/datasets/williamlucas0/employee-sample-data)

- **Function Definitions**

Using Any of the programming language implement below functions

    1. **createCollection(p_collection_name)**
    2. **indexData(p_collection_name, p_exclude_column)**:
       Index the given employee data into the specified collection, excluding the
       column provided in `p_exclude_column`.
    3. **searchByColumn(p_collection_name, p_column_name,
       p_column_value)**:
       Search within the specified collection for records where the column
       `p_column_name` matches the value `p_column_value`.
    4. **getEmpCount(p_collection_name)**
    5. **delEmpById(p_collection_name, p_employee_id)**
    6. · **getDepFacet(p_collection_name)**:
       Retrieve the count of employees grouped by department from the specified
       collection.

- **Function Executions**

Once the functions are implemented, execute the functions in the given order with the
parameters mentioned
    1. Var v_nameCollection = 'Hash_<Your Name>'
    2. Var v_phoneCollection ='Hash_<Your Phone last four digits'
    3. createCollection(v_nameCollection)
    4. createCollection(v_phoneCollection)
    5. getEmpCount(v_nameCollection)
    6. indexData**(v_**nameCollection,'Department')
    7. indexData**(v_** phoneCollection, 'Gender')
    8. delEmpById (v_ nameCollection ,'E02003')
    9. getEmpCount(v_nameCollection)
    10. searchByColumn(v_nameCollection,'Department','IT')
    11. searchByColumn(v_nameCollection,'Gender' ,'Male')
    12. searchByColumn(v_ phoneCollection,'Department','IT')
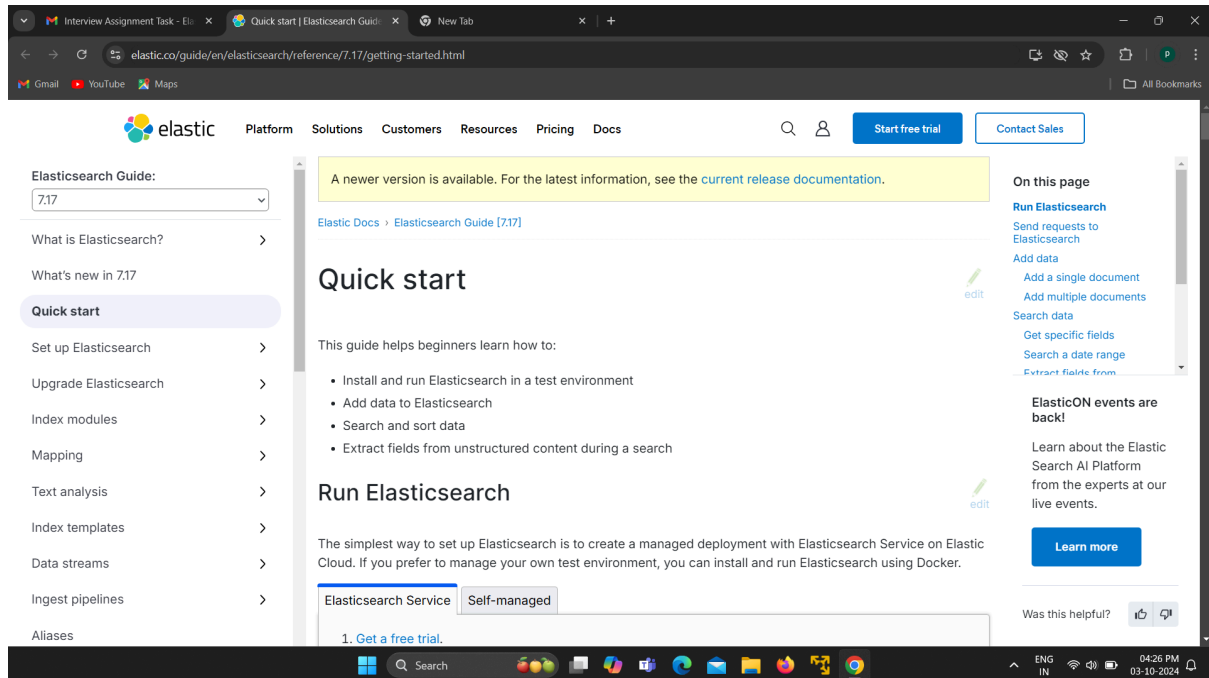    13. getDepFacet(v_ nameCollection)
    14. getDepFacet(v_ phoneCollection)

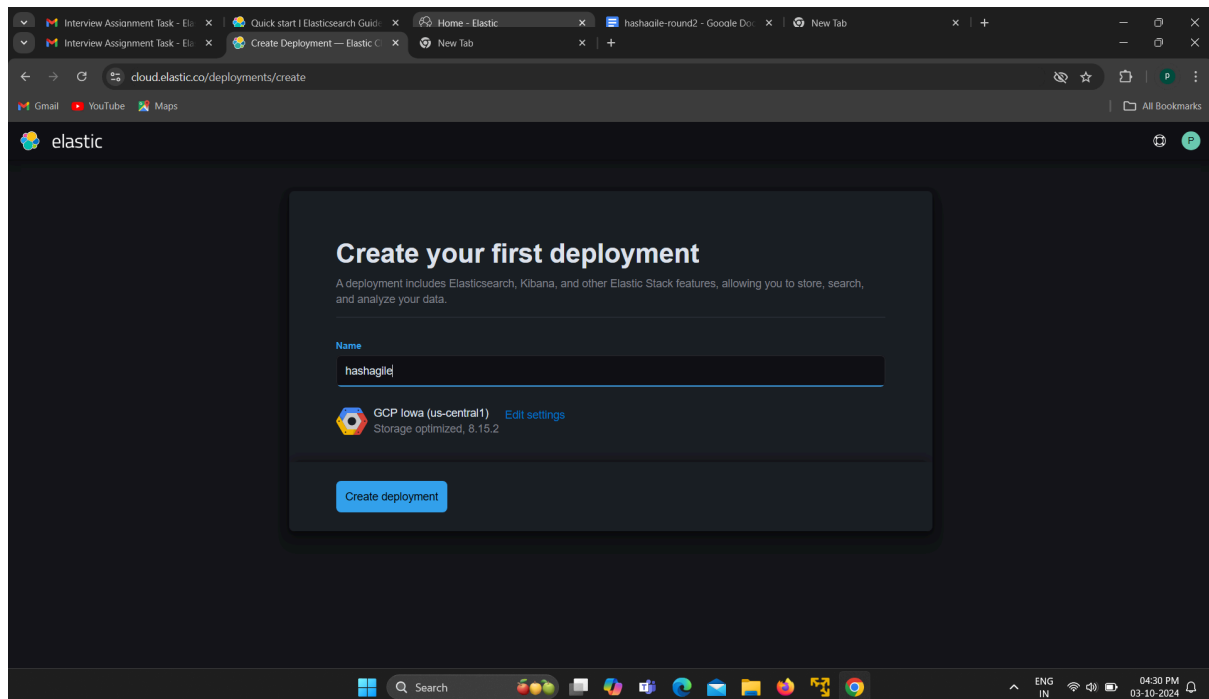# Assignment - Result(step-by-step):

- **Installation Instructions**
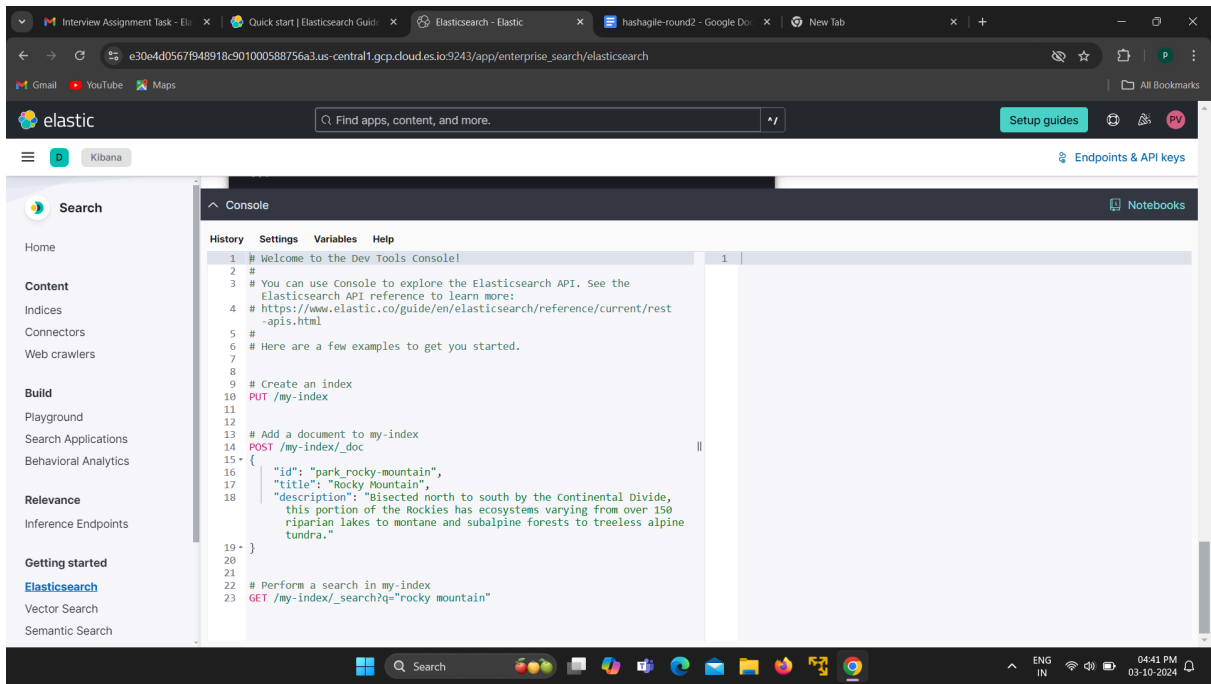    1. Read about Elasticsearch: [Getting Started with Elasticsearch](#).

2. Install Elasticsearch on your local machine.
3. Download Employee Data set from
   https://www.kaggle.com/datasets/williamlucas0/employee-sample-data
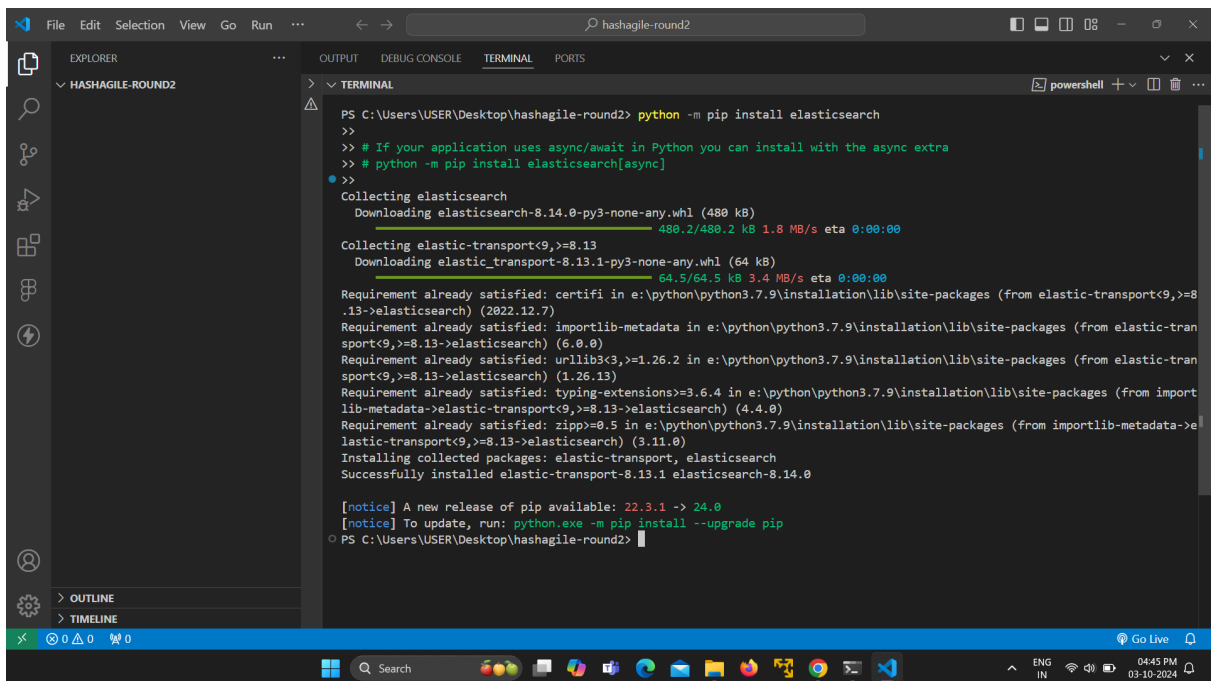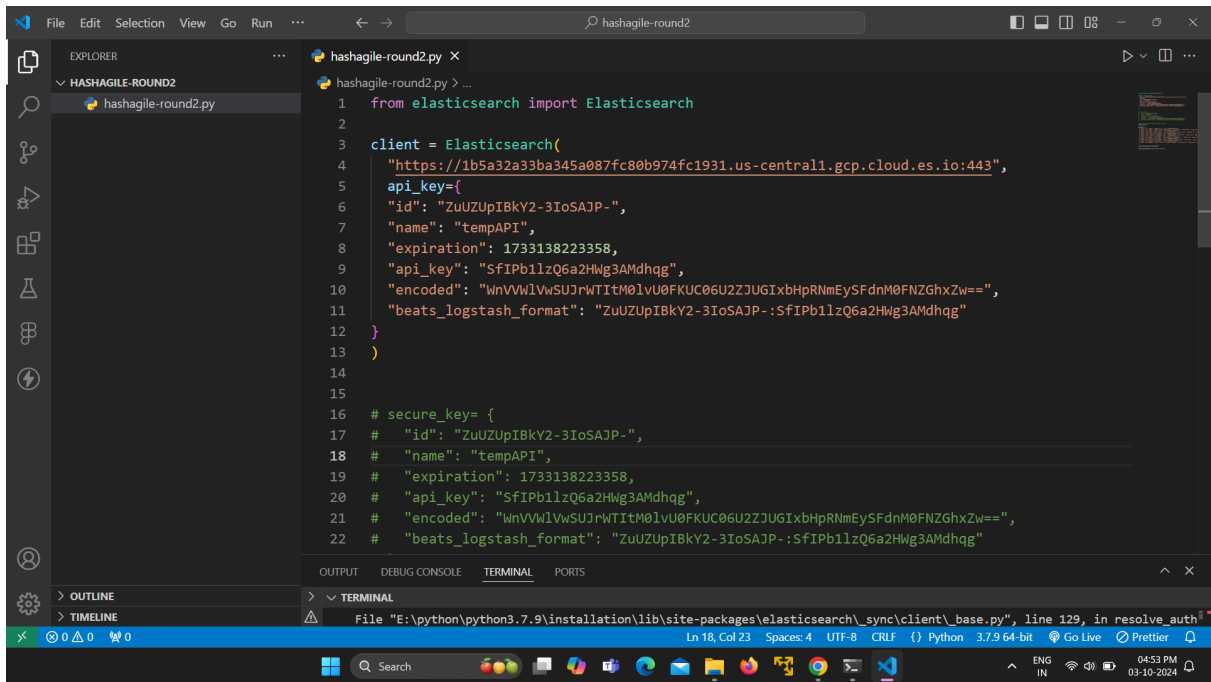
# Getting started with ElasticSearch:



# Making use of Elastic cloud(Alternative to local installation)

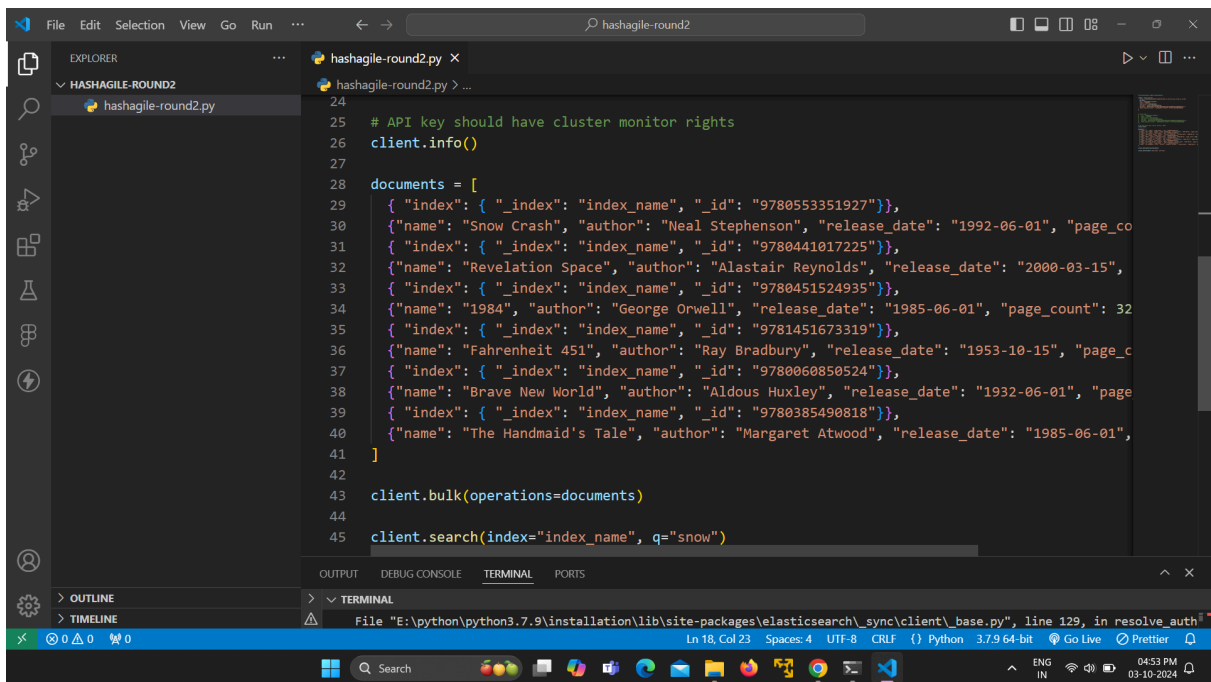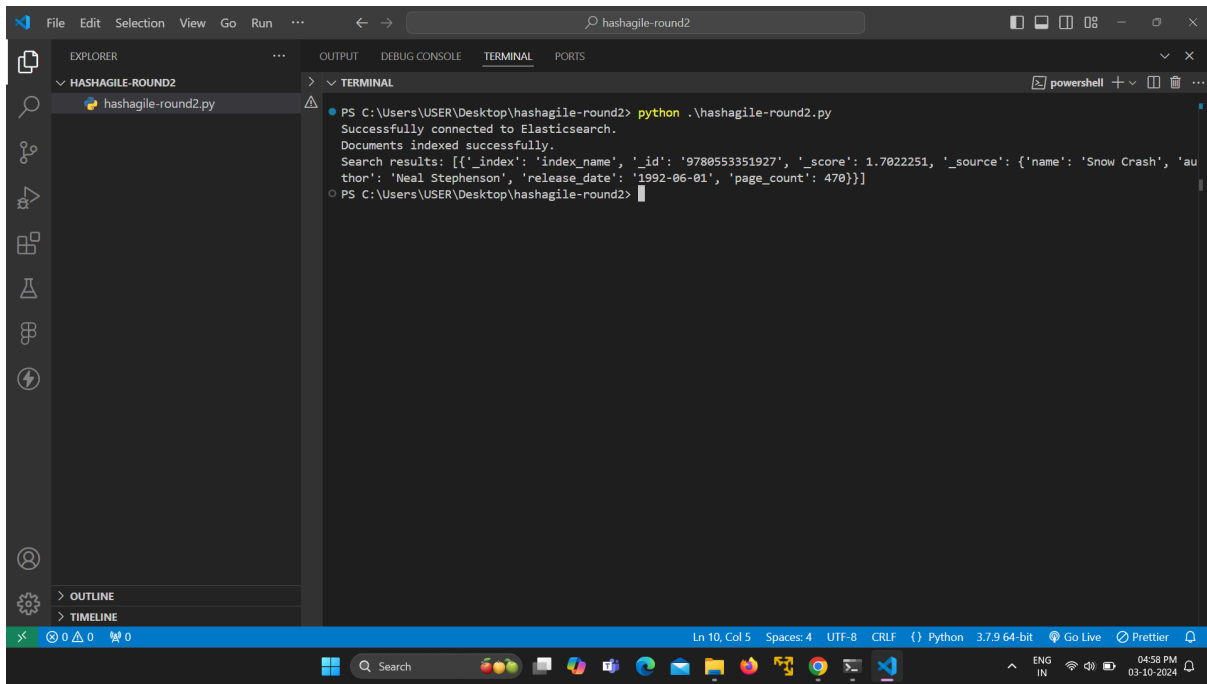Using this default code i am trying to stimulate in my local machine

```python
from elasticsearch import Elasticsearch

client = Elasticsearch(
    "https://1b5a32a33ba345a087fc80b974fc1931.us-central1.gcp.cloud.es.io:443",
    api_key={
    "id": "ZuUZUpIBkY2-3IoSAJP-",
    "name": "tempAPI",
    "expiration": 1733138223358,
    "api_key": "SfIPb1lzQ6a2HWg3AMdhqg",
    "encoded": "WnVVWlVwSUJrWTItM0lvU0FKUC06U2ZJUGIxbHpRNmEySFdnM0FNZGhxZw==",
    "beats_logstash_format": "ZuUZUpIBkY2-3IoSAJP-:SfIPb1lzQ6a2HWg3AMdhqg"
}
)


# secure_key= {
#    "id": "ZuUZUpIBkY2-3IoSAJP-",
#    "name": "tempAPI",
#    "expiration": 1733138223358,
#    "api_key": "SfIPb1lzQ6a2HWg3AMdhqg",
#    "encoded": "WnVVWlVwSUJrWTItM0lvU0FKUC06U2ZJUGIxbHpRNmEySFdnM0FNZGhxZw==",
#    "beats_logstash_format": "ZuUZUpIBkY2-3IoSAJP-:SfIPb1lzQ6a2HWg3AMdhqg"
```
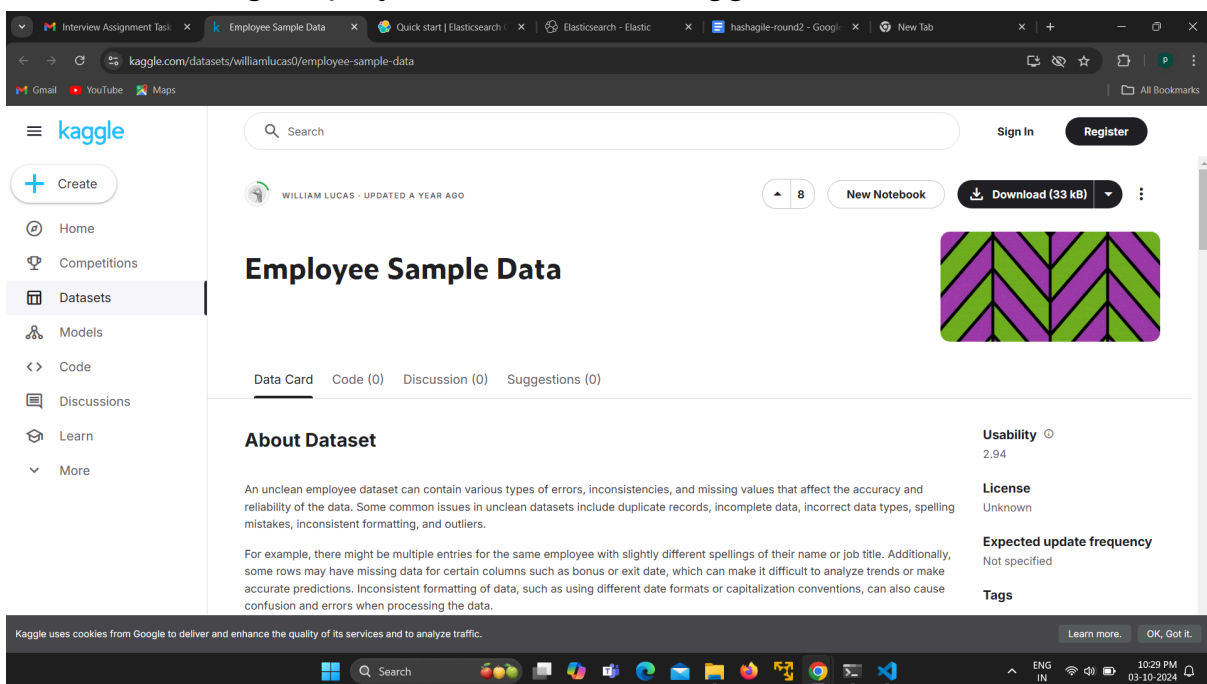
```
File "E:\python\python3.7.9\installation\lib\site-packages\elasticsearch\_sync\client\_base.py", line 129, in resolve_auth
```

```python
# API key should have cluster monitor rights
client.info()

documents = [
    { "index": { "_index": "index_name", "_id": "9780553351927"}},
    {"name": "Snow Crash", "author": "Neal Stephenson", "release_date": "1992-06-01", "page_co
    { "index": { "_index": "index_name", "_id": "9780441017225"}},
    {"name": "Revelation Space", "author": "Alastair Reynolds", "release_date": "2000-03-15",
    { "index": { "_index": "index_name", "_id": "9780451524935"}},
    {"name": "1984", "author": "George Orwell", "release_date": "1985-06-01", "page_count": 32
    { "index": { "_index": "index_name", "_id": "9781451673319"}},
    {"name": "Fahrenheit 451", "author": "Ray Bradbury", "release_date": "1953-10-15", "page_c
    { "index": { "_index": "index_name", "_id": "9780060850524"}},
    {"name": "Brave New World", "author": "Aldous Huxley", "release_date": "1932-06-01", "page
    { "index": { "_index": "index_name", "_id": "9780385490818"}},
    {"name": "The Handmaid's Tale", "author": "Margaret Atwood", "release_date": "1985-06-01",
]

client.bulk(operations=documents)

client.search(index="index_name", q="snow")
```

```
File "E:\python\python3.7.9\installation\lib\site-packages\elasticsearch\_sync\client\_base.py", line 129, in resolve_auth
```

# C. Downloading employee dataset from kaggle



Extract csv file from archive.zip and place in hashagile folder along with code file

# Function Definitions:

1. **createCollection(p_collection_name):**

```python
def create_collection(collection_name):
    if not client.indices.exists(index=collection_name):
        client.indices.create(index=collection_name)
        print(f"Collection '{collection_name}' created.")
    else:
        print(f"Collection '{collection_name}' already exists.")
```

2. **indexData(p_collection_name, p_exclude_column):**

```python
def index_data(collection_name, exclude_column):
    # Load the employee data
    try:
        df = pd.read_csv('employee_data.csv', encoding='ISO-8859-1')  # Adjust
encoding as needed
    except Exception as e:
        print(f"Error reading CSV file: {e}")
        return

    df = df.drop(columns=[exclude_column])  # Exclude the specified column
    df.fillna('', inplace=True)  # Replace NaNs with empty strings

    # Prepare documents for indexing
    documents = df.to_dict(orient='records')
    print(f"Preparing to index {len(documents)} documents in '{collection_name}'.")

    # Index data
    actions = []
    for i, doc in enumerate(documents):
        action = {
            '_op_type': 'index',
            '_index': collection_name,
            '_id': str(i),
            '_source': doc
        }
        actions.append(action)

    # Bulk index documents
    try:
        helpers.bulk(client, actions)
        print(f"Data indexed in '{collection_name}' excluding column
'{exclude_column}'.")
```

```
        except helpers.BulkIndexError as e:
            print(f"Error indexing documents: {e.errors}")
```

Index the given employee data into the specified collection, excluding the column provided in `p_exclude_column`.

3. **searchByColumn(p_collection_name, p_column_name, p_column_value)**:
   Search within the specified collection for records where the column `p_column_name` matches the value `p_column_value`.

```
def search_by_column(collection_name, column_name, column_value):
    query = {
        "query": {
            "match": {
                column_name: column_value
            }
        }
    }
    results = client.search(index=collection_name, body=query)
    return results['hits']['hits']
```

4. **getEmpCount(p_collection_name):**

```
def get_emp_count(collection_name):
    count = client.count(index=collection_name)
    return count['count']
```

5. **delEmpById(p_collection_name, p_employee_id):**

```
def del_emp_by_id(collection_name, employee_id):
    try:
        client.delete(index=collection_name, id=employee_id)
        print(f"Employee with ID '{employee_id}' deleted from '{collection_name}'.")
    except Exception as e:
        print(f"Error deleting employee: {e}")
```

6. · **getDepFacet(p_collection_name)**:

```
def get_dep_facet(collection_name):
    query = {
        "size": 0,
        "aggs": {
            "departments": {
                "terms": {
                    "field": "Department.keyword"
                }
            }
        }
    }
```

```
result = client.search(index=collection_name, body=query)
return result['aggregations']['departments']['buckets']
```
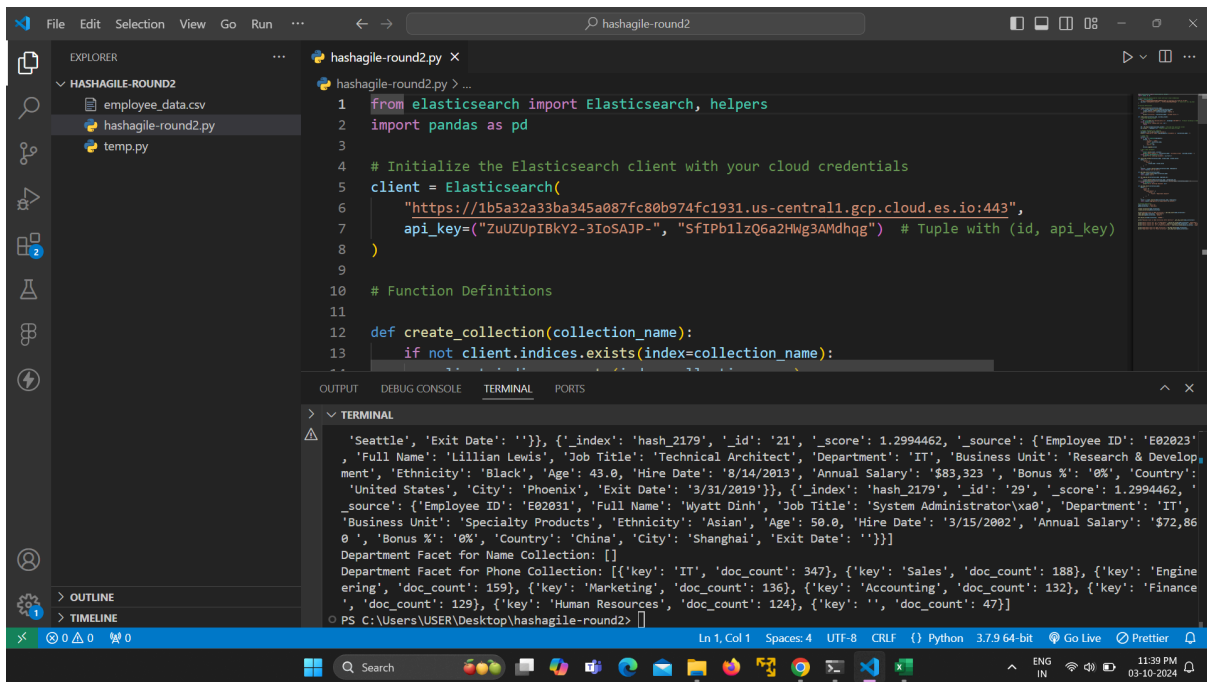
# Function Execution:

1. Var v_nameCollection = 'Hash_<Your Name>'
   name_collection = 'hash_ram'

2. Var v_phoneCollection ='Hash_<Your Phone last four digits'
   phone_collection = 'hash_2179'

3. createCollection(v_nameCollection)
   create_collection(name_collection)

4. createCollection(v_phoneCollection)
   create_collection(phone_collection)

5. getEmpCount(v_nameCollection)
   print("Employee Count in Name Collection:",
   get_emp_count(name_collection))

6. indexData(v_nameCollection,'Department')
   index_data(name_collection, 'Department')

7. indexData(v_ phoneCollection, 'Gender')
   index_data(phone_collection, 'Gender')

8. delEmpById (v_ nameCollection ,'E02003')
   del_emp_by_id(name_collection, 'E02003')

9. getEmpCount(v_nameCollection)
   print("Employee Count in Name Collection after deletion:",
   get_emp_count(name_collection))

10. searchByColumn(v_nameCollection,'Department','IT')
    print("Search results for 'IT' in Department:",
    search_by_column(name_collection, 'Department', 'IT'))

11. searchByColumn(v_nameCollection,'Gender' ,'Male')
    print("Search results for 'Male' in Gender:",
    search_by_column(name_collection, 'Gender', 'Male'))

12. searchByColumn(v_ phoneCollection,'Department','IT')
    print("Search results for 'IT' in Phone Collection:",
    search_by_column(phone_collection, 'Department', 'IT'))

13. getDepFacet(v_ nameCollection)

print("Department Facet for Name Collection:",
get_dep_facet(name_collection))

14. getDepFacet(v_phoneCollection)
print("Department Facet for Phone Collection:",
get_dep_facet(phone_collection))

Result: