# Dig Data Management Assignment 5

信科四 郑元嘉 1800920541

In this report, we're about to carry out three identical tasks on the same data `student.csv` via three various databases, `Redis` , `MongoDB` and `Neo4j` .
To get it straight for our quantity analysis of consumed time, we declare the three tasks here :

1. Insertion task: insert all data into the empty corresponding databases
2. Update task: update all data with all `sex` -> `'F'` and all `Mjob` -> `'teacher'`
3. Search task: find all original data where `sex` == `'F'` , `paid` == `'no'` and `internet` == `'yes'` , and then return all matched results

# Task1 - Redis

- **Insertion Task**
  (Partial source code)

```
1  def insert_bulk(file_path, redis_host, redis_port, redis_db_index):
2      global timestamp
3      client = redis.Redis(host=redis_host, port=redis_port, db=redis_db_inde
4
5      data = pd.read_csv(file_path, delimiter=';')
6      print('Load data {} with size of {}.'.format(file_path, data.shape[0]))
7
8      timestamp.stamp('insert')
9
10     for i, row in data.iterrows():
11         mapping = dict(row)
12         del mapping['id']
13         client.hmset(row['id'], mapping)
14
15     timestamp.stamp('insert')
```

By the means of Redis, we can only insert one datum each time by function `hmset()` , which means hardly could we run a 'bulk' task.
Here's some point worth attention. The csv file `student.csv` applies an abnormal delimiter `;` instead of `,` , so we need to set parameter `delimiter` of function `read_csv()` .

**Consumed Time :**

```
File/hadoopHw/hw5  master ✗
▶ python3 redis_manipulation.py -m insert
Load data ./student.csv with size of 1044.
Total time: 0.5204873085021973
Insertion time: 0.513737678527832
```

- **Update Task**

  (Partial source code)

```python
1   def update_bulk(file_path, redis_host, redis_port, redis_db_index):
2       global timestamp
3       client = redis.Redis(host=redis_host, port=redis_port, db=redis_db_inde
4
5       data = pd.read_csv(file_path, delimiter=';')
6       print('Load data {} with size of {}.'.format(file_path, data.size))
7
8       timestamp.stamp('update')
9
10      new_job, new_sex = 'teacher', 'F'
11      for i, row in data.iterrows():
12          mapping = {"Mjob": new_job, "sex": new_sex}
13          client.hmset(row['id'], mapping)
14
15      timestamp.stamp('update')
```

  Similar to insertion task, we can only update single one datum each time merely by function
  `hmset()` .

  **Consumed Time :**

```
File/hadoopHw/hw5  master ✗
▶ python3 redis_manipulation.py -m update
Load data ./student.csv with size of 35496.
Total time: 0.16183137893676758
Update time: 0.15511655807495117
```

- **Search Task**

  (Partial source code)

```python
1   def search_bulk(file_path, redis_host, redis_port, redis_db_index):
2       global timestamp
3       client = redis.Redis(host=redis_host, port=redis_port, db=redis_db_ind
4
5       timestamp.stamp('search')
6       sex_filter, paid_filter, internet_filter = 'F', 'no', 'yes'
7       output = []
8
9       keys = client.keys('*')
```

```
10        for key in keys:
11            datum = client.hgetall(key)
12            if datum['sex'] == sex_filter and datum['paid'] == paid_filter and
13                output.append(datum)
14
15        print('Search task returns {} matched items.'.format(len(output)))
16
17        timestamp.stamp('search')
```

Search task on Redis is quite unfriendly and inelegant. We query all keys first, and then query the data according to the keys, and ultimately we filter every datum by comparing the values of selected fields. On the whole, the filtering task is run on the client side which would potentially cause unnecessary and redundant transmission of data. Yet it can release the burden of Redis. It could be a trade-off.

**Consumed Time :**

```
File/hadoopHw/hw5  master ✗
▶ python3 redis_manipulation.py -m search
Search task returns 346 matched items.
Total time: 0.23855113983154297
Search time: 0.23825883865356445
```

# Task2 - MongoDB
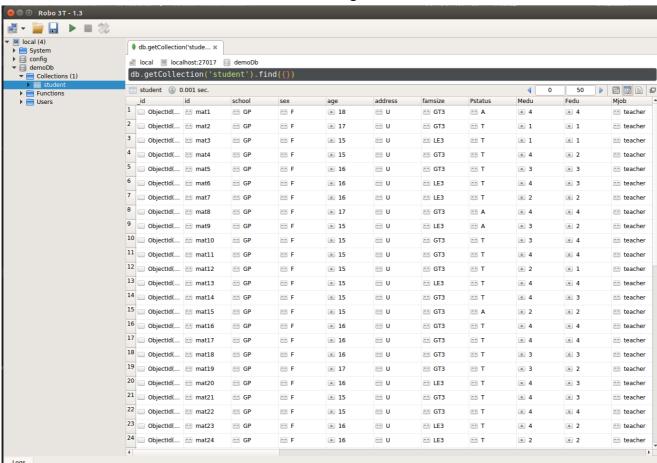
- **Insertion Task**
  (Partial source code)

```
1  def insert_bulk(file_path, db_host, db_port, db_name, db_col_name):
2      client = MongoClient(db_host, db_port)
3
4      db = client[db_name]
5      collection = db[db_col_name]
6
7      data = pd.read_csv(file_path, delimiter=';')
8      print('Load data {} with size of {}.'.format(file_path, data.shape[0]))
9
10     timestamp.stamp('insert')
11
12     bulk = []
13     for i, row in data.iterrows():
14         datum = dict(row)
15         bulk.append(datum)
16
```

```
17        collection.insert_many(bulk)
18
19        timestamp.stamp('insert')
```

First of all, we create one collection named `student` in our database. Unlike Redis, MongoDB supports bulk operation which reduce the waste of transmission time in a leap. We read the file `student.csv` , convert them into `dictionary` type, and call function `insert_many()` to be done with this task with ease.

Plus, we can use `Robot 3T` – a GUI tool for MongoDB – to have a look at the data.



**Consumed Time :**



- **Update Task**

  (Partial source code)

```
1  def update_bulk(file_path, db_host, db_port, db_name, db_col_name):
2      client = MongoClient(db_host, db_port)
3
4      db = client[db_name]
5      collection = db[db_col_name]
6
```

```
  7        global timestamp
  8        timestamp.stamp('update')
  9

 10        new_job, new_sex = 'teacher', 'F'
 11        res = collection.update_many({}, {'$set': {'Mjob': new_job, 'sex': nev
 12

 13        timestamp.stamp('update')
```

The story is almost the same. MongoDB supports bulk operation of update (modification).
Thus we're allowed to call function `update_many()` as soon as prepared for our conversion
preprocess. Because the client side only sends the filter query and wanted modification, all
operations such as retrieving and filtering are undergone in MongoDB which lessens the
workload on client sides and makes transmission more efficient.

**Consumed Time :**

```
File/hadoopHw/hw5  master ✗
▶ python3 mongodb_manipulation.py -m update
Total time: 0.009916543960571289
Update time: 0.00903630256652832
```

- **Search Task**

  (Partial source code)

```
  1  def search_bulk(file_path, db_host, db_port, db_name, db_col_name):
  2      client = MongoClient(db_host, db_port)
  3

  4      db = client[db_name]
  5      collection = db[db_col_name]
  6

  7      global timestamp
  8      timestamp.stamp('search')
  9

 10      sex_filter, paid_filter, internet_filter = 'F', 'no', 'yes'
 11      cursor = collection.find({'sex': sex_filter, 'paid': paid_filter, 'inte
 12      output = list(cursor)
 13

 14      print('Search task returns {} matched items.'.format(len(output)))
 15

 16      timestamp.stamp('search')
```

Bulk search task can be carried out via mere function `find()` supported by MongoDB. It's
noteworthy that the return of the function `find()` is a class `Cursor` rather than the actual
data, so we need to call extra functions triggering conversion like `list()` to retrieve the
actual data if needed. I dig into the source code of `pymongo` packages and find that `Cursor`
contains no actual data in RAM while when we call extra functions like `list()` on `Cursor`

instance, `pymongo` would send another request to MongoDB to retrieve the actual ones which is a clever and time-saving design in most cases.

**Consumed Time :**

```
File/hadoopHw/hw5  master ✗
▶ python3 mongodb_manipulation.py -m search
Search task returns 633 matched items.
Total time: 0.007836580276489258
Search time: 0.0066602230072021484
```

# Task3 - Neo4j

- **Insertion Task**

  (Partial source code)

```
 1  def insert_bulk(file_path, uri, acc, pwd):
 2      driver = GraphDatabase.driver(uri, auth=(acc, pwd))
 3
 4      data = pd.read_csv(file_path, delimiter=';')
 5      print('Load data {} with size of {}.'.format(file_path, data.shape[0]))
 6
 7      global timestamp
 8      timestamp.stamp('insert')
 9
10      statement = 'CREATE (:student{{{}}})'.format(','.join(['{v}: {{{v}}}'.f
11
12      with driver.session() as session:
13          for i, row in data.iterrows():
14              session.run(statement, dict(row))
15
16      timestamp.stamp('insert')
```

All manipulation of Neo4j is gone through scripts on Python. Due to lack of bulk insertion, we need to iteratively run the `CREATE` operation to insert all data. Plus, the rules of script format are kind of unfriendly and flexible which causes redundant programming codes for Python than other databases APIs.

Of course, you may attempt to write all `CREATE` clauses in one script to run which would reduce the times of transmission. The partial source code is like

```
 1  def insert_bulk(file_path, uri, acc, pwd):
 2      driver = GraphDatabase.driver(uri, auth=(acc, pwd))
 3
 4      data = pd.read_csv(file_path, delimiter=';')
```

```
 5        print('Load data {} with size of {}.'.format(file_path, data.shape[0])

 6

 7        global timestamp
 8        timestamp.stamp('insert')

 9

10        with driver.session() as session:
11            s = ''
12            for i, row in data.iterrows():
13                # print(dict(row))
14                statement = 'CREATE (:student{{{}}}) '.format(','.join(['{v}:
15                s += statement
16            session.run(s)

17

18        timestamp.stamp('insert')
```

However, much to our surprise, the consumed time increases badly somehow. Here's a guess - one script line would create only one process in Neoj4 to finish all clauses which means the process is overwhelmed by a great amount of memory and operation sources.

**Consumed Time :**

```
File/hadoopHw/hw5   master ✗
▶ python3 neo4j_manipulation.py -m insert
Load data ./student.csv with size of 1044.
Total time: 8.839545249938965
Insertion time: 8.79400086402893
```

- **Update Task**

(Partial source code)

```
 1  def update_bulk(file_path, uri, acc, pwd):
 2      driver = GraphDatabase.driver(uri, auth=(acc, pwd))
 3
 4      global timestamp
 5      timestamp.stamp('update')
 6
 7      new_job, new_sex = 'teacher', 'F'
 8      statement = 'MATCH (s:student{}) SET s.Mjob = {Mjob} SET s.sex = {sex}
 9
10      with driver.session() as session:
11          res = session.run(statement, Mjob=new_job, sex=new_sex)
12          data = list(res.records())
13
14      timestamp.stamp('update')
```

Neo4j supports bulk update in some extent. The only information in the script is our matching conditions and the new values which is highly similar to MongoDB.
Note that the resulting return is not a friendly object we're familiar with, thus we can access

the values of the returned data via some tricks.

**Consumed Time :**

```
File/hadoopHw/hw5  master ✗
▶ python3 neo4j_manipulation.py -m update
Total time: 0.33045244216918945
Update time: 0.2942204475402832
```

- **Search Task**

- (Partial source code)

```python
def search_bulk(file_path, uri, acc, pwd):
    driver = GraphDatabase.driver(uri, auth=(acc, pwd))

    global timestamp
    timestamp.stamp('search')

    sex_filter, paid_filter, internet_filter = 'F', 'no', 'yes'
    statement = 'MATCH (s:student{sex: {sex}, paid: {paid}, internet: {inte

    with driver.session() as session:
        res = session.run(statement, sex=sex_filter, paid=paid_filter, inte
        data = list(res.records())

        print('Search task returns {} matched items.'.format(len(data)))

    timestamp.stamp('search')
```

Neo4j provides `MATCH` clause to run our search task as well. We call `MATCH` clause and `RETURN` clause to retrieve the results. The call function `.records()` on the returned result to access the matched data.

**Consumed Time :**

```
File/hadoopHw/hw5  master ✗
▶ python3 neo4j_manipulation.py -m search
Search task returns 346 matched items.
Total time: 0.16910862922668457
Search time: 0.13193345069885254
```

# Consumed Time Analysis

|  | Insertion Task | Update Task | Search Task |
|---|---|---|---|
| Redis | 0.5137 | 0.1551 | 0.2383 |

|  | Insertion Task | Update Task | Search Task |
|---|---|---|---|
| MongoDB | **0.2869** | **0.0090** | **0.0067** |
| Neo4j | 8.7940 | 0.2942 | 0.1320 |

- Insertion Task
  The rare support of bulk insertion mainly contributes to the prominent performance of MongoDB. The simplified structure of storage helps Redis surpass Neo4j while the more complicated design of structure enables Neo4j to manipulate a rather more complex searching (matching) task.
- Update Task
  The cost of update is quite similar to the one of insertion, so the ranking is the same as the one for insertion task. It's worth attention that the consumed time of Neo4j shrinks the most significantly which shows the overhead of creating one node is a heavy workload.
- Search Task
  The relative ranking between Redis and Neo4j switches owing to lack of bulk research. Plus, the filtering task is run on the client side in the case of Redis, so we need to retrieve all data back and scan through every single one which is such an inefficiency in transmission.

Last but not least, keep in mind that we only go through extremely limited experiments to analyze which gives a huge bias towards MongoDB and doesn't give a guarantee of good performance in any cases.