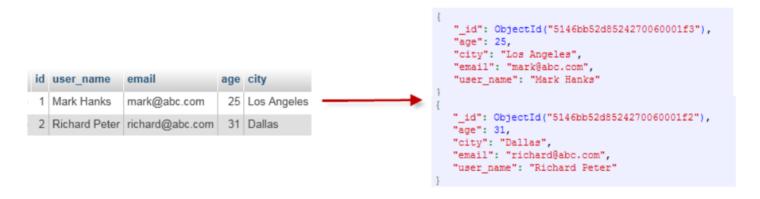
MongoDB使用

1.MongoDB概念

SQL 术语/概念	MongoDB 术语/概念	解释/说明
database	database	数据库
table	collection	数据库表/集合
row	document	数据记录行/文档
column	field	数据字段/域
index	index	索引
table joins		表连接,MongoDB不支持
primary key	primary key	主键,MongoDB自动将_id字段设置为主键

通过下图实例, 我们也可以更直观的了解Mongo中的一些概念:



2.MongoDB使用

1.安装相关库

```
1 pip install pymongo
2 # 清华镜像库
3 pip install pymongo -i https://pypi.tuna.tsinghua.edu.cn/simple --trusted-host
pypi.tuna.tsinghua.edu.cn
```

2.连接MongoDB

```
1 import pymongo
2
3 # 本地连接
4 url = "mongodb://localhost:27017/"
5 # 连接到本地 MongoDB 实例,指定数据库
6 url = "mongodb://localhost/mydatabase"
7 # 使用用户名和密码连接到本地 MongoDB 实例
8 url = "mongodb://username:password@localhost/mydatabase"
9 # 连接到远程 MongoDB 实例
10 url = "mongodb://remotehost:27017"
11 # 使用 SSL 连接到 MongoDB
12 url = "mongodb://username:password@localhost:27017/mydatabase?ssl=true"
13 # 使用多个选项连接
14 url = "mongodb://username:password@localhost:27017/mydatabase?
   authSource=admin&ssl=true"
15
16 myclient = pymongo.MongoClient(url)
```

3.数据库相关操作

1.创建数据库

在 MongoDB 中,数据库只有在内容插入后才会创建! 就是说,数据库创建后要创建集合(数据表)并插入一个文档(记录),数据库才会真正创建。

```
1 mydb = myclient["runoobdb"] # ""内为数据库名字
2
3 # 判断数据库是存在
4 dblist = myclient.list_database_names() # 获取所有的数据库名字
5 if "runoobdb" in dblist:
6 print("数据库已存在!")
```

4.集合(表)相关操作

1.创建

```
1 mydb = myclient["runoobdb"]
2 mycol = mydb["sites"]
3
```

```
4 # 判断集合是否存在
5 collist = mydb.list_collection_names()
6 if "sites" in collist:
7 print("集合已存在!")
```

2.写入数据

1.插入

insert_one()插入单条数据

```
1 # 插入单条数据
2 mydict = { "name": "RUNOOB", "alexa": "10000", "url": "https://www.runoob.com" }
3 x = mycol.insert_one(mydict)
4 # insert_one() 方法返回 InsertOneResult 对象,该对象包含 inserted_id 属性,它是插入文档的 id 值
5 print(x.inserted_id)
6 # 如果我们在插入文档时没有指定 _id, MongoDB 会为每个文档添加一个唯一的 id
```

_id	name	alexa	url
▶ 668a42cde4b836d1802fc7a8	RUNOOB	10000	https://www.runoob.com
668a430d8d93e863aa73275a	Google	1	https://www.google.com

insert_many()插入多条数据

```
1 mylist = [
 2 { "name": "Taobao", "alexa": "100", "url": "https://www.taobao.com" },
    { "name": "QQ", "alexa": "101", "url": "https://www.qq.com" },
 4 { "name": "Facebook", "alexa": "10", "url": "https://www.facebook.com" },
    { "name": "知乎", "alexa": "103", "url": "https://www.zhihu.com" },
    { "name": "Github", "alexa": "109", "url": "https://www.github.com" }
 6
7 ]
8 # 指定id插入
9 mylist = [
10 { "_id": 1, "name": "RUNOOB", "cn_name": "菜鸟教程"},
11 { "_id": 2, "name": "Google", "address": "Google 搜索"},
12 { "_id": 3, "name": "Facebook", "address": "脸书"},
13 { "_id": 4, "name": "Taobao", "address": "淘宝"},
    { "_id": 5, "name": "Zhihu", "address": "知乎"}
14
15 ]
16 # 插入多条数据
17 x = mycol.insert_many(mylist)
```

id	name	alexa	url	cn name	address
- 668a42cde4b836d1802fc7a8	RUNOOB	10000	https://www.runoob.com	(N/A)	(N/A)
668a430d8d93e863aa73275a	Google	1	https://www.google.com	(N/A)	(N/A)
668aa175af6d7a96bac77997	Taobao	100	https://www.taobao.com	(N/A)	(N/A)
668aa175af6d7a96bac77998	QQ	101	https://www.qq.com	(N/A)	(N/A)
668aa175af6d7a96bac77999	Facebook	10	https://www.facebook.com	(N/A)	(N/A)
668aa175af6d7a96bac7799a	知乎	103	https://www.zhihu.com	(N/A)	(N/A)
668aa175af6d7a96bac7799b	Github	109	https://www.github.com	(N/A)	(N/A)
1	RUNOOB	(N/A)	(N/A)	菜鸟教程	(N/A)
2	Google	(N/A)	(N/A)	(N/A)	Google 搜索
3	Facebook	(N/A)	(N/A)	(N/A)	脸书
4	Taobao	(N/A)	(N/A)	(N/A)	淘宝
5	Zhihu	(N/A)	(N/A)	(N/A)	知乎

2.更新

update_one()更新指定条件下的第一条数据

其中的参数myquery为通用的<mark>查询筛选器</mark>,用于查找到需要的文档

```
1 myquery = { "alexa": "10000" }
2 newvalues = { "$set": { "alexa": "12345" } }
3 # 为了证实只更改一个数据,将其中一个数据的alexa字段修改为10000
4 mycol.update_one(myquery, newvalues)
```

	_id	name	alexa	url
	668a42cde4b836d1802fc7a8	RUNOOB	10000	https://www.
	668a430d8d93e863aa73275a	Google	10000	https://www.
	668aa175af6d7a96bac77997	Taobao	100	https://www.
	668aa175af6d7a96bac77998	QQ	101	https://www.
Þ	668aa175af6d7a96bac77999	Facebook	10	https://www.
	668aa175af6d7a96bac7799a	知乎	103	https://www.
	668aa175af6d7a96bac7799b	Github	109	https://www.

_id	name	alexa	url
668a42cde4b836d1802fc7a8	RUNOOB	12345	https://v
668a430d8d93e863aa73275a	Google	10000	https://v
668aa175af6d7a96bac77997	Taobao	100	https://v
668aa175af6d7a96bac77998	QQ	101	https://v
668aa175af6d7a96bac77999	Facebook	10	https://
668aa175af6d7a96bac7799a	知乎	103	https://
668aa175af6d7a96bac7799b	Github	109	https://\

update_many()更新指定条件下的所有数据

```
1 myquery = { "name": { "$regex": "^G" } } # 正则表达式匹配开头为G的字段
2 newvalues = { "$set": { "alexa": "123" } }
3
4 x = mycol.update_many(myquery, newvalues)
```

-	-			
_id		name	alexa	ur
668a42cde4b836d1802fc7a8		RUNOOB	12345	ht
668a430d8d93e863aa73275a		Google	12345	ht
668aa175af6d7a96bac77997		Taobao	100	ht
668aa175af6d7a96bac77998		QQ	101	ht
668aa175af6d7a96bac77999		Facebook	10	ht
668aa175af6d7a96bac7799a		知乎	103	ht
668aa175af6d7a96bac7799b		Github	12345	ht
	1	RUNOOB	(N/A)	(N
	2	Google	12345	(N
	2	Facabaak	/N1/A3	/NI

3.替换

replace_one()替换指定条件下的第一条数据

□ 八知事为 ■ 又个 、 \/ 加尔 ↑= UE	יוי וויף וויי ידו)天尘颜已	で 会人 い 寺山	Ⅲ 刀切 ₩ 蚁流土川	6
_id	name	alexa	url		cn_name	address
668a42cde4b836d1802fc7a8	RUNOOB	12345	https://ww	w.runoob.com	(N/A)	(N/A)
▶ 668a430d8d93e863aa73275a	Xhs	111	https://ww	w.xiaohongshu.co	m (N/A)	(N/A)
668aa175af6d7a96bac77997	Taobao	100	https://ww	w.taobao.com	(N/A)	(N/A)
668aa175af6d7a96bac77998	QQ	101	https://ww	w.qq.com	(N/A)	(N/A)

4.删除

delete_one()删除查询条件下单条数据 delete_many()删除查询条件下所有数据 这里就不做演示了

```
1 myquery = { "name": "Xhs" }
2
```

```
3 mycol.delete_one(myquery)
4 mycol.delete_many(myquery)
```

5.批量操作

可以对某个集合使用bulk_write来合并使用上述操作

- insert_one()
- insert_many()
- update_one()
- update_many()
- replace_one()
- delete_one()
- delete_many()

```
1 operations = [
       pymongo.InsertOne(
2
 3
           {
                "name": "快手",
 4
               "alexa": "1234",
 5
                "url": "https://www.kuaishou.com"
 6
7
           }
8
       ),
9
       pymongo.UpdateMany(
           { "name": "知乎" },
10
           { "$set": { "alexa": "101" }},
11
12
       ),
13
       pymongo.DeleteOne(
           { "name": "Xhs" }
14
       )
15
16 ]
17 results = mycol.bulk_write(operations)
```

3.读取数据

1.查找函数

find_one()查找一条

find()查找全部

```
1 print(mycol.find_one())
2
3 for x in mycol.find():
4  print(x)
```

```
{'_id': ObjectId('668a42cde4b836d1802fc7a8'), 'name': 'RUNOOB', 'alexa': '12345', 'url': 'https://www.runoob.com'}

{'_id': ObjectId('668a42cde4b836d1802fc7a8'), 'name': 'RUNOOB', 'alexa': '12345', 'url': 'https://www.runoob.com'}

{'_id': ObjectId('668a430d8d93e863aa73275a'), 'name': 'Xhs', 'alexa': '111', 'url': 'https://www.xiaohongshu.com'}

{'_id': ObjectId('668aa175af6d7a96bac77997'), 'name': 'Taobao', 'alexa': '100', 'url': 'https://www.taobao.com'}

{'_id': ObjectId('668aa175af6d7a96bac77998'), 'name': 'QQ', 'alexa': '101', 'url': 'https://www.qq.com'}

{'_id': ObjectId('668aa175af6d7a96bac77999'), 'name': 'Facebook', 'alexa': '103', 'url': 'https://www.facebook.com'}

{'_id': ObjectId('668aa175af6d7a96bac7799a'), 'name': '知乎', 'alexa': '103', 'url': 'https://www.github.com'}

{'_id': ObjectId('668aa175af6d7a96bac7799b'), 'name': 'Github', 'alexa': '12345', 'url': 'https://www.github.com'}

{'_id': ObjectId('668aa175af6d7a96bac7799b'), 'name': 'Github', 'alexa': '12345', 'url': 'https://www.github.com'}

{'_id': ObjectId('668aa175af6d7a96bac7799b'), 'name': 'Github', 'alexa': '12345', 'url': 'https://www.github.com'}

{'_id': ObjectId('668aa175af6d7a96bac7799b'), 'name': 'IRUNOOB', 'alexa': '12345',
```

2.指定查询-查询筛选器

精确匹配

指定某个字段的某个值

```
1 for x in mycol.find({"name" : "Google"}):
2 # 多个字段在最外侧的{}内保证"key" : value,的格式即可
3 # for x in mycol.find({"name": "Google", "url": "111"})
4 print(x)
```

比较操作符

• \$gt: 大于

• \$lte: 小于或等于

• \$ne: 不等于

```
1 for x in mycol.find({"alexa": {"$ne" : "12345"}}):
2    print(x)
```

```
{'_id': ObjectId('668a430d8d93e863aa73275a'), 'name': 'Xhs', 'alexa': '111', 'url': 'https://www.xiaohongshu.com'}
{'_id': ObjectId('668aa175af6d7a96bac77997'), 'name': 'Taobao', 'alexa': '100', 'url': 'https://www.taobao.com'}
{'_id': ObjectId('668aa175af6d7a96bac77998'), 'name': 'QQ', 'alexa': '101', 'url': 'https://www.facebook.com'}
{'_id': ObjectId('668aa175af6d7a96bac77999'), 'name': 'Facebook', 'alexa': '10', 'url': 'https://www.facebook.com'}
{'_id': ObjectId('668aa175af6d7a96bac7799a'), 'name': '知乎', 'alexa': '103', 'url': 'https://www.zhihu.com'}
{'_id': 1, 'name': 'RUNOOB', 'cn_name': '菜鸟教程'}
{'_id': 3, 'name': 'Facebook', 'address': '脸书'}
{'_id': 4, 'name': 'Taobao', 'address': '知乎'}
```

逻辑操作符

- \$and ,返回符合*所有*子句条件的所有文档
- \$or ,返回符合一个子句条件的所有文档
- \$nor ,返回所有*不*符合任何子句条件的文档
- \$not ,返回与表达式*不*匹配的所有文档

```
{'_id': ObjectId('668a430d8d93e863aa73275a'), 'name': 'Xhs', 'alexa': '111', 'url': 'https://www.xiaohongshu.com'}
{'_id': ObjectId('668aa175af6d7a96bac77997'), 'name': 'Taobao', 'alexa': '100', 'url': 'https://www.taobao.com'}
{'_id': ObjectId('668aa175af6d7a96bac77998'), 'name': 'QQ', 'alexa': '101', 'url': 'https://www.qq.com'}
{'_id': ObjectId('668aa175af6d7a96bac77999'), 'name': 'Facebook', 'alexa': '10', 'url': 'https://www.facebook.com'}
{'_id': ObjectId('668aa175af6d7a96bac7799a'), 'name': '知乎', 'alexa': '103', 'url': 'https://www.zhihu.com'}
{'_id': 1, 'name': 'RUNOOB', 'cn_name': '菜鸟教程'}
{'_id': 3, 'name': 'Facebook', 'address': '脸书'}
{'_id': 4, 'name': 'Taobao', 'address': '淘宝'}
{'_id': 5, 'name': 'Zhihu', 'address': '知乎'}
```

数组操作符

- \$all ,它会返回带有数组的文档,而该数组包含查询中的所有元素
- \$elemMatch ,如果数组字段中的元素与查询中的所有条件匹配,则返回文档
- \$size ,返回包含指定大小数组的所有文档

元素操作符

- \$exists ,匹配具有指定字段的文档
- \$type ,如果字段为指定类型,则选择文档

评估操作符

求值操作符根据对单个字段或整个文集文件的求值结果返回数据

- \$text ,对文档执行文本Atlas Search
- \$regex ,返回与指定正则表达式匹配的文档
- \$mod ,它对字段的值执行模运算并返回余数为指定值的文档

3.指定要返回的字段

指定要包含的字段,默认情况包含 id字段

```
1 # { "<Field Name>": 1 }
2 for x in mycol.find({}, { "name" : 1, "url" : 1}):
3  print(x)
```

```
PS <u>D:\study\grade6\shixi\social</u>> python -u "d:\study\grade6\shixi\social\demo.py"
{'_id': ObjectId('668a42cde4b836d1802fc7a8'), 'name': 'RUNOOB', 'url': 'https://www.runoob.com'}
{'_id': ObjectId('668a430d8d93e863aa73275a'), 'name': 'Xhs', 'url': 'https://www.xiaohongshu.com'}
{'_id': ObjectId('668aa175af6d7a96bac77997'), 'name': 'Taobao', 'url': 'https://www.taobao.com'}
{'_id': ObjectId('668aa175af6d7a96bac77998'), 'name': 'QQ', 'url': 'https://www.qq.com'}
{'_id': ObjectId('668aa175af6d7a96bac77999'), 'name': 'Facebook', 'url': 'https://www.facebook.com'}
{'_id': ObjectId('668aa175af6d7a96bac7799a'), 'name': '知乎', 'url': 'https://www.zhihu.com'}
{'_id': ObjectId('668aa175af6d7a96bac7799b'), 'name': 'Github', 'url': 'https://www.github.com'}
```

指定要排除的字段

```
1 # { "<Field Name>": 0 }
2 for x in mycol.find({}, { "url" : 0}):
3 print(x)
```

除了_id,你不能在一个对象中同时指定0和1,如果你设置了一个字段为0,则其他都为1,反之亦然。

```
1 for x in mycol.find({}, { "_id" : 0, "name" : 1, "url" : 1}):
2 print(x)
3 # 下面这种写法会报错
4 # for x in mycol.find({}, { "name" : 0, "url" : 1}):
```

```
{'name': 'RUNOOB', 'url': 'https://www.runoob.com'}
{'name': 'Xhs', 'url': 'https://www.xiaohongshu.com'}
{'name': 'Taobao', 'url': 'https://www.taobao.com'}
{'name': 'QQ', 'url': 'https://www.qq.com'}
{'name': 'Facebook', 'url': 'https://www.facebook.com'}
{'name': '知乎', 'url': 'https://www.zhihu.com'}
{'name': 'Github', 'url': 'https://www.github.com'}
```

4.指定要返回的文档

读取操作返回的最大文档数limit

```
1 # 第一种
2 for x in mycol.find().limit(3):
3    print(x)
4 # 第二种
5 for x in mycol.find(limit=3):
6    print(x)
```

文档排序sort

需要保证是同种类型数据(如汉字和英文排序会报错)

```
1 # 默认升序
2 for x in mycol.find().sort("alexa"):
3     print(x)
4     # 降序
5 for x in mycol.find().sort("alexa", pymongo.DESCENDING):
6     print(x)
7     # 第二种方式
8 for x in mycol.find(sort={"alexa", pymongo.ASCENDING}):
9     print(x)
```

先跳过某个数量文档,再查询skip

```
1 # 跳过前两个文档
2 for x in mycol.find().skip(2):
3     print(x)
4 # 第二种写法
5 for x in mycol.find({}, skip=2):
6     print(x)
```

这些方法可以组合,但执行顺序为排序,跳过,限制

```
1 # 第一种写法
2 for x in mycol.find().skip(2).sort("alexa", pymongo.DESCENDING).limit(3):
3    print(x)
4 # 第二种写法
5 for x in mycol.find(skip=2, sort={"alexa", pymongo.DESCENDING}, limit=3):
6    print(x)
```

_id	name	alexa	url	cn_name	address
668a42cde4b836d1802fc7a8	RUNOOB	12345	https://www.runoob.com	(N/A)	(N/A)
668a430d8d93e863aa73275a	Xhs	111	https://www.xiaohongshu.com	(N/A)	(N/A)
668aa175af6d7a96bac77997	Taobao	100	https://www.taobao.com	(N/A)	(N/A)
668aa175af6d7a96bac77998	QQ	101	https://www.qq.com	(N/A)	(N/A)
668aa175af6d7a96bac77999	Facebook	10	https://www.facebook.com	(N/A)	(N/A)
668aa175af6d7a96bac7799a	知乎	103	https://www.zhihu.com	(N/A)	(N/A)
668aa175af6d7a96bac7799b	Github	12345	https://www.github.com	(N/A)	(N/A)
	RUNOOB	(N/A)	(N/A)	菜鸟教程	(N/A)
2	2 Google	12345	(N/A)	(N/A)	Google 搜索
3	B Facebook	(N/A)	(N/A)	(N/A)	脸书
4	1 Taobao	(N/A)	(N/A)	(N/A)	淘宝
í	Zhihu	(N/A)	(N/A)	(N/A)	知乎

```
{'_id': 2, 'name': 'Google', 'address': 'Google 搜索', 'alexa': '12345'}
{'_id': ObjectId('668a430d8d93e863aa73275a'), 'name': 'Xhs', 'alexa': '111', 'url': 'https://www.xiaohongshu.com'}
{'_id': ObjectId('668aa175af6d7a96bac7799a'), 'name': '知乎', 'alexa': '103', 'url': 'https://www.zhihu.com'}
```

5.多表联合查询

一对一的表连接

即A表与B表连接的属性为B表的主键,即A中连接的值在B表中唯一

以以下两个集合为例,用A表orders的product_id和B表products的id作为连接,来获取每个order包含的具体产品信息

	_id	customer_id	orderdate	product_id	value
١	668fb8e96e7cdd1d91a3dbc7	elise_smith@myemail.com	2020-05-30 08:35:52.000	a1b2c3d4	431.43
	668fb8e96e7cdd1d91a3dbc8	tj@wheresmyemail.com	2019-05-28 19:13:32.000	z9y8x7w6	5.01
	668fb8e96e7cdd1d91a3dbc9	oranieri@warmmail.com	2020-01-01 08:25:37.000	ff11gg22hh33	63.13
	668fb8e96e7cdd1d91a3dbca	jjones@tepidmail.com	2020-12-26 08:55:46.000	a1b2c3d4	429.65

	_id	id	name	category	description
١	668fb8e96e7cdd1d91a3dbcb	a1b2c3d4	Asus Laptop	ELECTRONICS	Good value laptop for st
	668fb8e96e7cdd1d91a3dbcc	z9y8x7w6	The Day Of The Triffids	BOOKS	Classic post-apocalyptic
	668fb8e96e7cdd1d91a3dbcd	ff11gg22hh33	Morphy Richardds Food	KITCHENWARE	Luxury mixer turning go
	668fb8e96e7cdd1d91a3dbce	pqr678st	Karcher Hose Set	GARDEN	Hose + nosels + winder

```
1 orders_coll = mydb["orders"] # 集合A
 2 products_coll = mydb["products"] # 集合B
 3
 4 pipeline = []
 5
 6 pipeline.append({
 7
      "orderdate": {
 8
 9
              "$gte": datetime(2020, 1, 1, 0, 0, 0),
              "$lt": datetime(2021, 1, 1, 0, 0, 0)
10
11
          }
      }
12
13 })
14
15 pipeline.append({
      "$lookup": {
16
          "from": "products", # 需要联合的集合B
17
          "localField": "product_id", # 集合A的字段
18
          "foreignField": "id", # 集合B的字段
19
          "as": "product_mapping"
20
      }
21
22 })
23 # 等价于以下mysql语句
24 # SELECT *, product_mapping
25 # FROM orders
26 # WHERE product_mapping IN (SELECT *
27 # FROM products
28 # WHERE id= orders.product_id);
29 pipeline.extend([
      {
30
          # 将product_mapping字段设置为上一个$lookup阶段中创建的product_mapping对象中
31
   的第一个元素
          # 因为在B表中值唯一, 所以只取1个
32
33
          "$set": {
              "product_mapping": {"$first": "$product_mapping"}
34
35
          }
      },
36
37
          "$set": { # 根据product_mapping对象字段中的值创建两个新字段product_name和
38
   product_category
39
              "product_name": "$product_mapping.name",
              "product_category": "$product_mapping.category"
40
          }
41
      }
42
43 ])
44 # 从文档中删除不必要的字段
```

```
45 pipeline.append({"$unset": ["_id", "product_id", "product_mapping"]})
46
47 aggregation_result = orders_coll.aggregate(pipeline)
```

```
1 {
 2
     'customer_id': 'elise_smith@myemail.com',
     'orderdate': datetime.datetime(2020, 5, 30, 8, 35, 52),
 3
 4
     'value': 431.43,
     'product_name': 'Asus Laptop',
 5
     'product_category': 'ELECTRONICS'
 6
 7 }
 8 {
     'customer_id': 'oranieri@warmmail.com',
     'orderdate': datetime.datetime(2020, 1, 1, 8, 25, 37),
10
     'value': 63.13,
11
     'product_name': 'Morphy Richardds Food Mixer',
12
     'product_category': 'KITCHENWARE'
13
14 }
15 {
     'customer_id': 'jjones@tepidmail.com',
16
     'orderdate': datetime.datetime(2020, 12, 26, 8, 55, 46),
17
18
     'value': 429.65,
19
     'product_name': 'Asus Laptop',
     'product_category': 'ELECTRONICS'
20
21 }
```

多字段的表连接

以以下两个集合为例,用A表products的variation和name字段与B表orders的product_name和 product_variation连接,来获取某个product产生的所有订单

_id	variation	name	category	description
669092eb91b6f198e1c7ad51	Ultra HD	Asus Laptop	ELECTRONICS	Great for watching movi
669092eb91b6f198e1c7ad52	Standard Display	Asus Laptop	ELECTRONICS	Good value laptop for st
669092eb91b6f198e1c7ad53	1st Edition	The Day Of The Triffids	BOOKS	Classic post-apocalyptic
669092eb91b6f198e1c7ad54	2nd Edition	The Day Of The Triffids	BOOKS	Classic post-apocalyptic
669092eb91b6f198e1c7ad55	Deluxe	Morphy Richards Food 1	KITCHENWARE	Luxury mixer turning go

_id	customer_id	orderdate	product_name	product_variation	value
669092eb91b6f198e1c7ad56	elise_smith@myemail.com	2020-05-30 08:35:52.000	Asus Laptop	Standard Display	431.43
669092eb91b6f198e1c7ad57	tj@wheresmyemail.com	2019-05-28 19:13:32.000	The Day Of The Triffids	2nd Edition	5.01
669092eb91b6f198e1c7ad58	oranieri@warmmail.com	2020-01-01 08:25:37.000	Morphy Richards Food 1	Deluxe	63.13
669092eb91b6f198e1c7ad59	jjones@tepidmail.com	2020-12-26 08:55:46.000	Asus Laptop	Standard Display	429.65

```
1 pipeline = []
```

```
2 # 这里的embedded_pl是使用在lookup中,所以都是对要连接的B表进行处理
3 embedded_pl = [
       {
 4
 5
           "$match": { # 匹配连接的B表中的连接字段并重命名
              "$expr": {
 6
                  "$and": [
 7
 8
                      {"$eq": ["$product_name", "$$prdname"]},
                      {"$eq": ["$product_variation", "$$prdvartn"]}
9
10
                  ]
              }
11
12
          }
      }
13
14 ]
15
16 embedded_pl.append({
      "$match": { # 筛选B表中数据
17
          "orderdate": {
18
19
              "$gte": datetime(2020, 1, 1, 0, 0, 0),
              "$lt": datetime(2021, 1, 1, 0, 0, 0)
20
21
          }
22
      }
23 })
24
25 embedded_pl.append({ # 选择B表中不需要的字段,两表的连接字段会出现两次,所以需要删除一
   张表中的数据
     "$unset": ["_id", "product_name", "product_variation"]
26
27 })
28
29 pipeline.append({ # 这里是多字段的匹配处理
      "$lookup": {
30
          "from": "orders",
31
          "let": {
32
              "prdname": "$name",
33
              "prdvartn": "$variation"
34
35
          },
36
          "pipeline": embedded_pl,
          "as": "orders"
37
38
      }
39 })
40
41 pipeline.append({
       "$match": { # 对连接完后的数据进行筛选
42
          "orders": {"$ne": []}
43
44
      }
45 })
46
47 pipeline.append({ # 最后选择需要的字段
```

```
48    "$unset": ["_id", "description"]
49 })
50
51 aggregation_result = products_coll.aggregate(pipeline)
52
53 for document in aggregation_result:
54    print(document)
```

```
1 {
 2
     'name': 'Asus Laptop',
 3
     'variation': 'Standard Display',
 4
     'category': 'ELECTRONICS',
     'orders': [
 5
 6
 7
         'customer_id': 'elise_smith@myemail.com',
 8
         'orderdate': datetime.datetime(2020, 5, 30, 8, 35, 52),
         'value': 431.43
9
       },
10
11
         'customer_id': 'jjones@tepidmail.com',
12
         'orderdate': datetime.datetime(2020, 12, 26, 8, 55, 46),
13
         'value': 429.65
14
15
       }
     ٦
16
17 }
18 {
19
     'name': 'Morphy Richards Food Mixer',
20
     'variation': 'Deluxe',
     'category': 'KITCHENWARE',
21
22
     'orders': [
23
     {
         'customer_id': 'oranieri@warmmail.com',
24
         'orderdate': datetime.datetime(2020, 1, 1, 8, 25, 37),
25
         'value': 63.13
26
      }
27
    ]
28
29 }
```

5.GridFS-大数据(图片和视频)的上传与下载

如果文件大小超过16 MB 的 BSON 文档大小限制,则应使用 GridFS 存储桶包含以下集合,使用 GridFS 规范中定义的约定命名:

• chunks 集合存储二进制文件数据段。

• files 集合存储文件元数据。

id

使用 GridFS 存储文件时,驱动程序会将文件拆分成较小的数据块,每个数据块由 chunks 集合中的单独文档表示。它还在 files 集合中创建文档,其中包含文件 ID、文件名和其他文件元数据。

这里使用下方大小为28MB的视频作为示例

ZhihuSpider	2024/6/27 14:35	文件夹	
🕯 demo.py	2024/7/12 21:19	Python 源文件	1 KB
🕯 kshot.py	2024/7/3 22:40	Python 源文件	5 KB
ksVideoDownload.py	2024/6/26 10:29	Python 源文件	2 KB
MediaCrawler.zip	2024/7/1 13:05	ZIP 压缩文件	388,658 KB
ques.txt	2024/6/26 17:05	文本文档	1 KB
tempCodeRunnerFile.py	2024/7/10 19:25	Python 源文件	1 KB
video.mp4	2024/7/12 11:37	MP4 文件	29,684 KB
video1.mp4	2024/7/12 21:19	MP4 文件	29,684 KB
zhihucrawler.py	2024/6/27 10:30	Python 源文件	5 KB
動 项目架构 nna	2024/6/24 14:18	PNG 图片文件	119 KR



```
1 import pymongo
2 from gridfs import GridFS
3
4 myclient = pymongo.MongoClient('mongodb://localhost:27017/')
5
6 mydb = myclient["runoobdb"]
7 fs = GridFS(mydb) # 创建该数据库下的GridFS桶
8
9 myfile = open(file='D:/study/grade6/shixi/social/video.mp4', mode='rb') # 本地上 传
10 mydata = myfile.read()
11
12 myfile = fs.put(data=mydata, filename="video.mp4") # 这里的filename命名为 GridFS.files中的filename值
13
14 print(myfile) # 这里myfile返回的是_id值
```

66912d05479a7ce228a8edc8	video.mp4	261120	30396	6205 202	24-07-12 13:17:57.68
	III Tu a de			() 	
□ 开始事务 □ 文本 ▼ 〒 筛选 □ 排序 □ 付 ● 全部折叠 ● 类型颜色 □ 日 <p< td=""><td>₩ 分析 ₩ 数据生成</td></p<>				₩ 分析 ₩ 数据生成	
_id	files_id			n	data
▶ 66912d05479a7ce228a8edc9	66912d05479	a7ce228a8edc8		0	(BLOB) 255.00 KB
66912d05479a7ce228a8edca	66912d05479	a7ce228a8edc8		1	(BLOB) 255.00 KB
66912d05479a7ce228a8edcb	66912d05479	a7ce228a8edc8		2	(BLOB) 255.00 KB
66912d05479a7ce228a8edcc	66912d05479	a7ce228a8edc8		3	(BLOB) 255.00 KB
66912d05479a7ce228a8edcd	66912d05479	a7ce228a8edc8		4	(BLOB) 255.00 KB

filename

chunkSize

length

uploadDate

66912d05479a7ce228a8ee38	66912d05479a7ce228a8edc8	111	(BLOB) 255.00 KB
66912d05479a7ce228a8ee39	66912d05479a7ce228a8edc8	112	(BLOB) 255.00 KB
66912d05479a7ce228a8ee3a	66912d05479a7ce228a8edc8	113	(BLOB) 255.00 KB
66912d05479a7ce228a8ee3b	66912d05479a7ce228a8edc8	114	(BLOB) 255.00 KB
66912d05479a7ce228a8ee3c	66912d05479a7ce228a8edc8	115	(BLOB) 255.00 KB
▶ 66912d05479a7ce228a8ee3d	66912d05479a7ce228a8edc8	116	(BLOB) 103.79 KB

可以看到fs.files的_id字段值是和fs.chunks的filles_id是对应的,查找的时候也是通过该值来查找的

```
1 mydata = fs.get_version("video.mp4").read()
2 myfile = open(file='D:/study/grade6/shixi/social/video1.mp4', mode='wb')
3 rs = myfile.write(mydata)
```

i video.mp4	2024/7/12 11:37	MP4 文件	29,684 KB
■ video1.mp4	2024/7/12 21:19	MP4 文件	29,684 KB
- · · · · · · · · · · · · · · · · · · ·		- 1 25 10	