

# 1 Lab3 MongoDB 的聚合管道

## 1.1 何为聚合操作

聚合操作主要是通过对数据进行分组后做出一些简单的运算，例如平均，求和，最值等。MongoDB 中的聚合运算主要通过 `aggregate()` 方法实现

Python 中 `aggregate()` 方式实现了利用聚合管道对文档进行变化计算和展示。文档进入聚合管道会依次经过筛选 (filtering)，分组 (grouping) 并聚合，排序 (sorting)，投射 (projecting)，限制 (limiting) (或者跳过 (skipping) ) 变化。

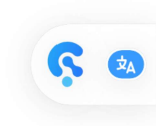
符号	含义	描述
\$match	筛选	按照一定条件筛选出特定的文档记录
\$project	选择	修改输入文档的结构。可以用来重命名、增加或删除域，也可以用于创建计算结果以及嵌套文档
\$group	分组	对文档按照字段分组，以便做一些聚合运算
\$sort	排序	对文档按照字段排序
\$limit	限制	限制MongoDB聚合管道返回的文档数
\$skip	跳过	在聚合管道中跳过指定数量的文档，并返回余下的文档

## 1.2 简单的聚合

首先连接数据库创建数据表并且插入数据

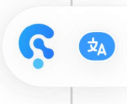
```
In [3]: # 连接Mongodb数据库
import pymongo

client = pymongo.MongoClient("mongodb://ecnu10215501412:ECNU10215501412@172.16.14.60")
db = client["ecnu10215501412"]
users_col = db["users"]
```



In [4]: # 插入数据

```
users_col.delete_many({})
users = [
    {
        "name": "Joe",
        "gender": "m",
        "age": 23,
        "birthdate": {"day": 15, "month": 3, "year": 1997},
        "hobby": ["football", "basketball", "reading"],
        "city": "Beijing",
        "time": [9, 18],
    },
    {
        "name": "Kate",
        "gender": "f",
        "age": 22,
        "birthdate": {"day": 25, "month": 7, "year": 1998},
        "hobby": ["reading", "piano"],
        "city": "Hangzhou",
        "time": [8, 17],
    },
    {
        "name": "Rose",
        "gender": "f",
        "age": 24,
        "birthdate": {"day": 3, "month": 3, "year": 1996},
        "hobby": ["basketball", "running", "traveling"],
        "city": "Shanghai",
        "time": [9, 19],
    },
    {
        "name": "Jason",
        "gender": "m",
        "age": 21,
        "birthdate": {"day": 17, "month": 12, "year": 1999},
        "hobby": ["cooking", "photography"],
        "city": "Chengdu",
        "time": [8, 20],
    },
    {
        "name": "Grace",
        "gender": "f",
        "age": 22,
        "birthdate": {"day": 10, "month": 6, "year": 1998},
        "hobby": ["photography", "cooking", "drama"],
        "city": "Nanjing",
        "time": [9, 18],
    },
    {
        "name": "Jessica",
        "gender": "f",
        "age": 22,
        "birthdate": {"day": 21, "month": 3, "year": 1998},
        "hobby": ["cooking", "piano"],
        "city": "Shanghai",
        "time": [10, 19],
    },
    {
        "name": "Donna",
        "gender": "f",
        "age": 22,
```



```

        "birthdate": {"day": 24, "month": 9, "year": 1998},
        "hobby": ["violin", "drama"],
        "city": "Shanghai",
        "time": [9, 20],
    },
    {
        "name": "Apple",
        "gender": "m",
        "age": 23,
        "birthdate": {"day": 20, "month": 9, "year": 1997},
        "hobby": ["violin", "running"],
        "city": "Chengdu",
        "time": [9, 19],
    },
    {
        "name": "Baba",
        "gender": "f",
        "age": 25,
        "birthdate": {"day": 20, "month": 9, "year": 1995},
        "hobby": ["violin", "basketball"],
        "city": "Chengdu",
        "time": [10, 19],
    },
]

users_col.insert_many(users)
content = users_col.find()
for each in content:
    print(each)

```

```

{'_id': ObjectId('651509dfbe9ac71da3f94365'), 'name': 'Joe', 'gender': 'm', 'age': 23, 'birthdate': {'day': 15, 'month': 3, 'year': 1997}, 'hobby': ['football', 'basketball', 'reading'], 'city': 'Beijing', 'time': [9, 18]}
{'_id': ObjectId('651509dfbe9ac71da3f94366'), 'name': 'Kate', 'gender': 'f', 'age': 22, 'birthdate': {'day': 25, 'month': 7, 'year': 1998}, 'hobby': ['reading', 'piano'], 'city': 'Hangzhou', 'time': [8, 17]}
{'_id': ObjectId('651509dfbe9ac71da3f94367'), 'name': 'Rose', 'gender': 'f', 'age': 24, 'birthdate': {'day': 3, 'month': 3, 'year': 1996}, 'hobby': ['basketball', 'running', 'traveling'], 'city': 'Shanghai', 'time': [9, 19]}
{'_id': ObjectId('651509dfbe9ac71da3f94368'), 'name': 'Jason', 'gender': 'm', 'age': 21, 'birthdate': {'day': 17, 'month': 12, 'year': 1999}, 'hobby': ['cooking', 'photography'], 'city': 'Chengdu', 'time': [8, 20]}
{'_id': ObjectId('651509dfbe9ac71da3f94369'), 'name': 'Grace', 'gender': 'f', 'age': 22, 'birthdate': {'day': 10, 'month': 6, 'year': 1998}, 'hobby': ['photography', 'cooking', 'drama'], 'city': 'Nanjing', 'time': [9, 18]}
{'_id': ObjectId('651509dfbe9ac71da3f9436a'), 'name': 'Jessica', 'gender': 'f', 'age': 22, 'birthdate': {'day': 21, 'month': 3, 'year': 1998}, 'hobby': ['cooking', 'piano'], 'city': 'Shanghai', 'time': [10, 19]}
{'_id': ObjectId('651509dfbe9ac71da3f9436b'), 'name': 'Donna', 'gender': 'f', 'age': 22, 'birthdate': {'day': 24, 'month': 9, 'year': 1998}, 'hobby': ['violin', 'drama'], 'city': 'Shanghai', 'time': [9, 20]}
{'_id': ObjectId('651509dfbe9ac71da3f9436c'), 'name': 'Apple', 'gender': 'm', 'age': 23, 'birthdate': {'day': 20, 'month': 9, 'year': 1997}, 'hobby': ['violin', 'running'], 'city': 'Chengdu', 'time': [9, 19]}
{'_id': ObjectId('651509dfbe9ac71da3f9436d'), 'name': 'Baba', 'gender': 'f', 'age': 25, 'birthdate': {'day': 20, 'month': 9, 'year': 1995}, 'hobby': ['violin', 'basketball'], 'city': 'Chengdu', 'time': [10, 19]}

```

按照城市分组计数

```
In [5]: result = users_col.aggregate([{"$group": {"_id": "$city", "count": {"$sum": 1}}}])
for each in result:
    print(each)

{'_id': 'Beijing', 'count': 1}
{'_id': 'Chengdu', 'count': 3}
{'_id': 'Nanjing', 'count': 1}
{'_id': 'Hangzhou', 'count': 1}
{'_id': 'Shanghai', 'count': 3}
```

从以上代码中，聚合管道中只有 group 一个操作。在 group 中，可以看到是按照 city 字段进行分组，最后通过“加一”聚合来实现分组计数的。在group中指定了两个字段，第一个是主键 '\_id'，来源于 city，第二个是 count，来源于求和。当然也可以根据我们的需要，修改/增减字段。

以下列出了常用的聚集运算

符号	含义
\$sum	求和
\$avg	求平均
\$min	最小值
\$max	最大值
\$push	聚合成数组
\$addToSet	聚合成几何
\$first	排序取第一个
\$last	排序取最后一个

### 1.2.1 练习

Task 1 计算不同性别用户的平均年龄，最大年龄，最小年龄并且输出

```
In [8]: # todo
result = users_col.aggregate([{"$group": {"_id": "$gender", "ave_age": {"$avg": "$age", "$max_age": "$max_age", "$min_age": "$min_age"}}}])
for each in result:
    print(each)
"""
目标结果
{'_id': 'f', 'avg_age': 22.833333333333332, 'max_age': 25, 'min_age': 22}
{'_id': 'm', 'avg_age': 22.333333333333332, 'max_age': 23, 'min_age': 21}
"""

{'_id': 'm', 'ave_age': 22.333333333333332, 'max_age': 23, 'min_age': 21}
{'_id': 'f', 'ave_age': 22.833333333333332, 'max_age': 25, 'min_age': 22}
```

Out[8]: "\n目标结果\n{'\_id': 'f', 'avg\_age': 22.833333333333332, 'max\_age': 25, 'min\_age': 22}\n{'\_id': 'm', 'avg\_age': 22.333333333333332, 'max\_age': 23, 'min\_age': 21}\n"

## Task 2 列出不同性别同学名单 (提示: \$push)

```
In [9]: # todo
result = users_col.aggregate([{"$group": {"_id": "$gender", "name": {"$push": "$name"}}}]
for each in result:
    print(each)
"""
目标结果
{'_id': 'm', 'list': ['Joe', 'Jason', 'Apple']}
{'_id': 'f', 'list': ['Kate', 'Rose', 'Grace', 'Jessica', 'Donna', 'Baba']}
"""

{'_id': 'm', 'name': ['Joe', 'Jason', 'Apple']}
{'_id': 'f', 'name': ['Kate', 'Rose', 'Grace', 'Jessica', 'Donna', 'Baba']}
```

```
Out[9]: "\n目标结果\n{'_id': 'm', 'list': ['Joe', 'Jason', 'Apple']}\n{'_id': 'f', 'list': ['Kate', 'Rose', 'Grace', 'Jessica', 'Donna', 'Baba']}\n"
```

## 1.3 复杂查询

一个完整的查询一般需要经过:

1. 通过条件筛选文档记录 (选择文档的行记录)
2. 分组并聚合
3. 对文档按照某些字段排序
4. 调整文档的键值对形式 (调整文档的列)
5. 通过limit或者skip展示特定数量的记录

例如: 筛选年龄大于等于20岁的同学, 并将这些学生的按照城市分组计算平均年龄后升序排列, 显示城市和平均年龄

```
In [10]: match = {"$match": {"age": {"$gte": 20}}}
group = {"$group": {"_id": "$city", "avg_age": {"$avg": "$age"}}}
sort = {"$sort": {"avg_age": 1}} # 1代表升序, -1代表降序
project = {"$project": {"avg_age": 1}}
result = users_col.aggregate([match, group, sort, project])
for each in result:
    print(each)

{'_id': 'Hangzhou', 'avg_age': 22.0}
{'_id': 'Nanjing', 'avg_age': 22.0}
{'_id': 'Shanghai', 'avg_age': 22.666666666666668}
{'_id': 'Chengdu', 'avg_age': 23.0}
{'_id': 'Beijing', 'avg_age': 23.0}
```

格式化输出每个同学的生日日期, 按照生日日期排序

```
In [11]: sort = {"$sort": {"birthdate.year": -1, "birthdate.month": -1, "birthdate.day": -1}}
project = {
    "$project": {
        "_id": 0,
        "name": 1,
        "birthday": {
            "$concat": [
                {"$toString": "$birthdate.year"},
                "-",
                {"$toString": "$birthdate.month"},
                "-",
                {"$toString": "$birthdate.day"}
            ]
        }
    }
} # 0代表不显示该字段, 1代表显示该字段
result = users_col.aggregate([sort, project])
for each in result:
    print(each)
```

```
{'name': 'Jason', 'birthday': '1999-12-12'}
{'name': 'Donna', 'birthday': '1998-9-9'}
{'name': 'Kate', 'birthday': '1998-7-7'}
{'name': 'Grace', 'birthday': '1998-6-6'}
{'name': 'Jessica', 'birthday': '1998-3-3'}
{'name': 'Apple', 'birthday': '1997-9-9'}
{'name': 'Joe', 'birthday': '1997-3-3'}
{'name': 'Rose', 'birthday': '1996-3-3'}
{'name': 'Baba', 'birthday': '1995-9-9'}
```

以上的'toString'的作用是将数字转换成字符串, 更多函数可以参见

[https://blog.csdn.net/weixin\\_43632687/article/details/104201185](https://blog.csdn.net/weixin_43632687/article/details/104201185)  
([https://blog.csdn.net/weixin\\_43632687/article/details/104201185](https://blog.csdn.net/weixin_43632687/article/details/104201185))

按照城市, 性别分组计数

```
In [12]: group = {
    "$group": {"_id": {"city": "$city", "gender": "$gender"}, "count": {"$sum": 1}}
}
project = {
    "$project": {
        "_id": 0,
        "city": "$_id.city",
        "gender": "$_id.gender",
        "count": "$count",
    }
}
result = users_col.aggregate([group, project])
for each in result:
    print(each)
```

```
{'city': 'Hangzhou', 'gender': 'f', 'count': 1}
{'city': 'Shanghai', 'gender': 'f', 'count': 3}
{'city': 'Beijing', 'gender': 'm', 'count': 1}
{'city': 'Chengdu', 'gender': 'm', 'count': 2}
{'city': 'Chengdu', 'gender': 'f', 'count': 1}
{'city': 'Nanjing', 'gender': 'f', 'count': 1}
```

以上用到的重命名方法和多字段分组聚合的方法需要好好体会

### 1.3.1 练习



Task 3 找出喜欢'violin'的人数 (提示: \$in)

```
In [28]: # todo
match = {"$match": {"hobby": {"$in": ['violin']}}}
group = {"$group": {"_id": "like_violin", "count": {"$sum": 1}}}

result = users_col.aggregate([match, group])
for each in result:
    print(each)
"""
目标结果
{'_id': 'like_violin', 'count': 3}
"""

{'_id': 'like_violin', 'count': 3}
```

Out[28]: "\n目标结果\n{'\_id': 'like\_violin', 'count': 3}\n"

\$unwind 拆分数组, 查询拥有各个爱好的学生人数

```
In [29]: # 先通过数组拆分将一条记录拆分成多条记录
unwind = {"$unwind": "$hobby"}
group = {"$group": {"_id": "$hobby", "count": {"$sum": 1}}}
result = users_col.aggregate([unwind, group])
for each in result:
    print(each)

{'_id': 'football', 'count': 1}
{'_id': 'violin', 'count': 3}
{'_id': 'basketball', 'count': 3}
{'_id': 'reading', 'count': 2}
{'_id': 'running', 'count': 2}
{'_id': 'cooking', 'count': 3}
{'_id': 'traveling', 'count': 1}
{'_id': 'photography', 'count': 2}
{'_id': 'drama', 'count': 2}
{'_id': 'piano', 'count': 2}
```

### 1.3.2 练习

Task 4 找出爱好个数为3的同学, 展示姓名, 年龄与爱好 (不使用 \$size 来求长度, 要求使用 \$unwind 来拆分数组和 \$push 来合并数组)



```
In [51]: # todo
unwind = {"$unwind": "$hobby"}
group={"$group":{"_id":{"name": "$name", "age": "$age"},'hobby':{'$push':"$hobby"}},'
project1 = {
    "$project": {
        "_id": 0,
        "name": "$name",
        "age": "$age",
        "hobby":1,
        "count": 1,
    }
}
match={"$match":{"count":3}}
project2 = {
    "$project": {
        "_id": 0,
        "name": "$_id.name",
        "age": "$_id.age",
        "hobby": "$hobby",
    }
}
result = users_col.aggregate([unwind, project1, group, match, project2])
for each in result:
    print(each)
"""
目标结果
{'name': 'Rose', 'age': 24, 'hobby': ['basketball', 'running', 'traveling']}
{'name': 'Grace', 'age': 22, 'hobby': ['photography', 'cooking', 'drama']}
{'name': 'Joe', 'age': 23, 'hobby': ['football', 'basketball', 'reading']}
"""
```

```
{'name': 'Rose', 'age': 24, 'hobby': ['basketball', 'running', 'traveling']}
{'name': 'Joe', 'age': 23, 'hobby': ['football', 'basketball', 'reading']}
{'name': 'Grace', 'age': 22, 'hobby': ['photography', 'cooking', 'drama']}
```

```
Out[51]: "\n目标结果\n{'name': 'Rose', 'age': 24, 'hobby': ['basketball', 'running', 'traveling']}\n{'name': 'Grace', 'age': 22, 'hobby': ['photography', 'cooking', 'drama']}\n{'name': 'Joe', 'age': 23, 'hobby': ['football', 'basketball', 'reading']}\n"
```

## 1.4 索引：对 Lab2 的一些补充

### 1.4.1 单条索引：按照年龄生序

```
In [52]: users_col.create_index([("name", 1)], unique=True)
```

```
Out[52]: 'name_1'
```

### 1.4.2 复合索引：创建复合索引，按照姓名生序，按照年龄降序

```
In [53]: users_col.create_index([("name", 1), ("age", -1)], unique=True)
```

```
Out[53]: 'name_1_age_-1'
```

### 1.4.3 删除索引

```
In [54]: print("删除前索引信息\n", users_col.index_information())
users_col.drop_index("name_1") # 括号里面的参数是索引名
users_col.drop_index("name_1_age_-1")
print("删除后索引信息\n", users_col.index_information())
```

删除前索引信息

```
{'_id': {'v': 2, 'key': [(' _id', 1)]}, 'name_1': {'v': 2, 'key': [(' name', 1)],
'unique': True}, 'name_1_age_-1': {'v': 2, 'key': [(' name', 1), (' age', -1)], 'un
ique': True}}
```

删除后索引信息

```
{'_id': {'v': 2, 'key': [(' _id', 1)]}}
```

### 1.4.4 性能测试

索引的价值在于提高基于索引字进行段数据查询的效率，我们可以通过构造一批数据，对比建索引前的查询时间来体会索引的价值

```
In [55]: # 构造3000000条数据
import random

# 先清空一下数据库
users_col.delete_many({})

batch_users = []
sex = ["f", "m"]
for i in range(300000):
    user = {
        "name": "xxx" + str(i),
        "age": random.randint(20, 55), # 产生20, 55之间的随机数
        "gender": sex[random.randint(0, 1)],
    }
    batch_users.append(user)
users_col.insert_many(batch_users)
```

Out[55]: <pymongo.results.InsertManyResult at 0x7f7a8d1f67a0>

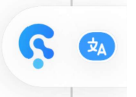
```
In [56]: # 直接查询用时
import datetime

starttime = datetime.datetime.now()

result = users_col.find(
    {
        "$or": [
            {"name": "xxx10000"},
            {"name": "xxx140000"},
            {"name": "xxx9000"},
            {"name": "xxx23000"},
            {"name": "xxx24050"},
            {"name": "xxx12000"},
            {"name": "xxx14300"},
            {"name": "xxx9300"},
            {"name": "xxx23300"},
            {"name": "xxx24350"},
            {"name": "xxx11100"},
            {"name": "xxx15200"},
            {"name": "xxx8100"},
            {"name": "xxx22100"},
            {"name": "xxx26150"},
            {"name": "xxx10200"},
            {"name": "xxx14020"},
            {"name": "xxx9020"},
            {"name": "xxx23020"},
            {"name": "xxx24070"},
            {"name": "xxx10300"},
            {"name": "xxx14030"},
            {"name": "xxx9030"},
            {"name": "xxx23030"},
            {"name": "xxx24080"},
        ]
    }
)

for each in result:
    print(each)
endtime = datetime.datetime.now()

print("开始时间:", starttime)
print("结束时间:", endtime)
print("时间差 (微秒):", (endtime - starttime).microseconds)
```

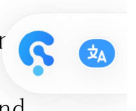


```
{ '_id': ObjectId('651517acbe9ac71da3f96312'), 'name': 'xxx8100', 'age': 40, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96696'), 'name': 'xxx9000', 'age': 22, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f966aa'), 'name': 'xxx9020', 'age': 25, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f966b4'), 'name': 'xxx9030', 'age': 21, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f967c2'), 'name': 'xxx9300', 'age': 37, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96a7e'), 'name': 'xxx10000', 'age': 54, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96b46'), 'name': 'xxx10200', 'age': 53, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f96baa'), 'name': 'xxx10300', 'age': 34, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96eca'), 'name': 'xxx11100', 'age': 27, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f9724e'), 'name': 'xxx12000', 'age': 50, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f97a32'), 'name': 'xxx14020', 'age': 51, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f97a3c'), 'name': 'xxx14030', 'age': 43, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f97b4a'), 'name': 'xxx14300', 'age': 30, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f97ece'), 'name': 'xxx15200', 'age': 35, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f999c2'), 'name': 'xxx22100', 'age': 48, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f99d46'), 'name': 'xxx23000', 'age': 28, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f99d5a'), 'name': 'xxx23020', 'age': 34, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f99d64'), 'name': 'xxx23030', 'age': 43, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f99e72'), 'name': 'xxx23300', 'age': 36, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f9a160'), 'name': 'xxx24050', 'age': 31, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a174'), 'name': 'xxx24070', 'age': 24, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a17e'), 'name': 'xxx24080', 'age': 49, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a28c'), 'name': 'xxx24350', 'age': 24, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a994'), 'name': 'xxx26150', 'age': 31, 'gender': 'f' }
{ '_id': ObjectId('651517adbe9ac71da3fb664e'), 'name': 'xxx140000', 'age': 52, 'gender': 'm' }
```

开始时间: 2023-09-28 06:05:36.745759

结束时间: 2023-09-28 06:05:36.906517

时间差(微秒): 160758



```
In [57]: # 创建索引查询用时间
users_col.create_index(["name", 1], unique=True)
starttime = datetime.datetime.now()
result = users_col.find(
    {
        "$or": [
            {"name": "xxx10000"},
            {"name": "xxx140000"},
            {"name": "xxx9000"},
            {"name": "xxx23000"},
            {"name": "xxx24050"},
            {"name": "xxx12000"},
            {"name": "xxx14300"},
            {"name": "xxx9300"},
            {"name": "xxx23300"},
            {"name": "xxx24350"},
            {"name": "xxx11100"},
            {"name": "xxx15200"},
            {"name": "xxx8100"},
            {"name": "xxx22100"},
            {"name": "xxx26150"},
            {"name": "xxx10200"},
            {"name": "xxx14020"},
            {"name": "xxx9020"},
            {"name": "xxx23020"},
            {"name": "xxx24070"},
            {"name": "xxx10300"},
            {"name": "xxx14030"},
            {"name": "xxx9030"},
            {"name": "xxx23030"},
            {"name": "xxx24080"},
        ]
    }
)
for each in result:
    print(each)
endtime = datetime.datetime.now()

print("开始时间:", starttime)
print("结束时间:", endtime)
print("时间差 (微秒):", (endtime - starttime).microseconds)
users_col.drop_index("name_1") # 结束后删除索引以防之后忘记删除
```



```
{ '_id': ObjectId('651517acbe9ac71da3f96a7e'), 'name': 'xxx10000', 'age': 54, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96b46'), 'name': 'xxx10200', 'age': 53, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f96baa'), 'name': 'xxx10300', 'age': 34, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96eca'), 'name': 'xxx11100', 'age': 27, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f9724e'), 'name': 'xxx12000', 'age': 50, 'gender': 'f' }
{ '_id': ObjectId('651517adbe9ac71da3fb664e'), 'name': 'xxx140000', 'age': 52, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f97a32'), 'name': 'xxx14020', 'age': 51, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f97a3c'), 'name': 'xxx14030', 'age': 43, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f97b4a'), 'name': 'xxx14300', 'age': 30, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f97ece'), 'name': 'xxx15200', 'age': 35, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f999c2'), 'name': 'xxx22100', 'age': 48, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f99d46'), 'name': 'xxx23000', 'age': 28, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f99d5a'), 'name': 'xxx23020', 'age': 34, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f99d64'), 'name': 'xxx23030', 'age': 43, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f99e72'), 'name': 'xxx23300', 'age': 36, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f9a160'), 'name': 'xxx24050', 'age': 31, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a174'), 'name': 'xxx24070', 'age': 24, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a17e'), 'name': 'xxx24080', 'age': 49, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a28c'), 'name': 'xxx24350', 'age': 24, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f9a994'), 'name': 'xxx26150', 'age': 31, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96312'), 'name': 'xxx8100', 'age': 40, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f96696'), 'name': 'xxx9000', 'age': 22, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f966aa'), 'name': 'xxx9020', 'age': 25, 'gender': 'm' }
{ '_id': ObjectId('651517acbe9ac71da3f966b4'), 'name': 'xxx9030', 'age': 21, 'gender': 'f' }
{ '_id': ObjectId('651517acbe9ac71da3f967c2'), 'name': 'xxx9300', 'age': 37, 'gender': 'f' }
```

开始时间: 2023-09-28 06:05:41.780469  
结束时间: 2023-09-28 06:05:41.789171  
时间差(微秒): 8702

在我的机器上运行后发现建立索引前后明显的查询时间分别为307881微秒和11248, 创建索引之后的查询耗时有了明显的降低

```
In [58]: users_col.delete_many({})
# 插入无关数据后记得删除，保持良好习惯，以免学院服务器崩坏
```

```
Out[58]: <pymongo.results.DeleteResult at 0x7f7aa520acb0>
```



## 1.4.5 练习

Task 5 下面需要同学们探索对比，创建索引对插入数据的影响

```
In [59]: # todo

import random
import datetime

# 先清空一下数据库
users_col.delete_many({})

starttime = datetime.datetime.now()
# 构造3000000条数据，计时
batch_users = []
sex = ["f", "m"]
for i in range(300000):
    user = {
        "name": "xxx" + str(i),
        "age": random.randint(20, 55), # 产生20, 55之间的随机数
        "gender": sex[random.randint(0, 1)],
    }
    batch_users.append(user)
users_col.insert_many(batch_users)

endtime = datetime.datetime.now()

print("开始时间:", starttime)
print("结束时间:", endtime)
print("不建立索引时间差（微秒）:", (endtime - starttime).microseconds)

users_col.delete_many({})
```

开始时间: 2023-09-28 06:18:01.616151

结束时间: 2023-09-28 06:18:06.530770

不建立索引时间差（微秒）: 914619

```
Out[59]: <pymongo.results.DeleteResult at 0x7f7a6fc6ef20>
```

```
In [60]: # todo

import random
import datetime

# 先清空一下数据库
users_col.delete_many({})

# 建立索引
users_col.create_index([("name", 1)], unique=True)

starttime = datetime.datetime.now()
# 构造3000000条数据, 计时
batch_users = []
sex = ["f", "m"]
for i in range(300000):
    user = {
        "name": "xxx" + str(i),
        "age": random.randint(20, 55), # 产生20, 55之间的随机数
        "gender": sex[random.randint(0, 1)],
    }
    batch_users.append(user)
users_col.insert_many(batch_users)

endtime = datetime.datetime.now()

print("开始时间:", starttime)
print("结束时间:", endtime)
print("建立索引时间差(微秒):", (endtime - starttime).microseconds)
users_col.drop_index("name_1") # 结束后删除索引以防之后忘记删除

users_col.delete_many({})
"""
过程&结论
"""
```

```
开始时间: 2023-09-28 06:18:15.672993
结束时间: 2023-09-28 06:18:21.801404
建立索引时间差(微秒): 128411
```

```
Out[60]: '\n过程&结论\n'
```

上述两段代码的运行结果分别为:

- 开始时间: 2023-09-28 06:18:01.616151
- 结束时间: 2023-09-28 06:18:06.530770
- 不建立索引时间差(微秒): 914619
- 开始时间: 2023-09-28 06:18:15.672993
- 结束时间: 2023-09-28 06:18:21.801404
- 建立索引时间差(微秒): 128411

因此, 创建索引之后的插入耗时有了明显的降低



