

به نام خدا



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



درس کلان داده

تمرین شماره دو

نام و نام خانوادگی : ?

شماره دانشجویی : ?

خرداد 99

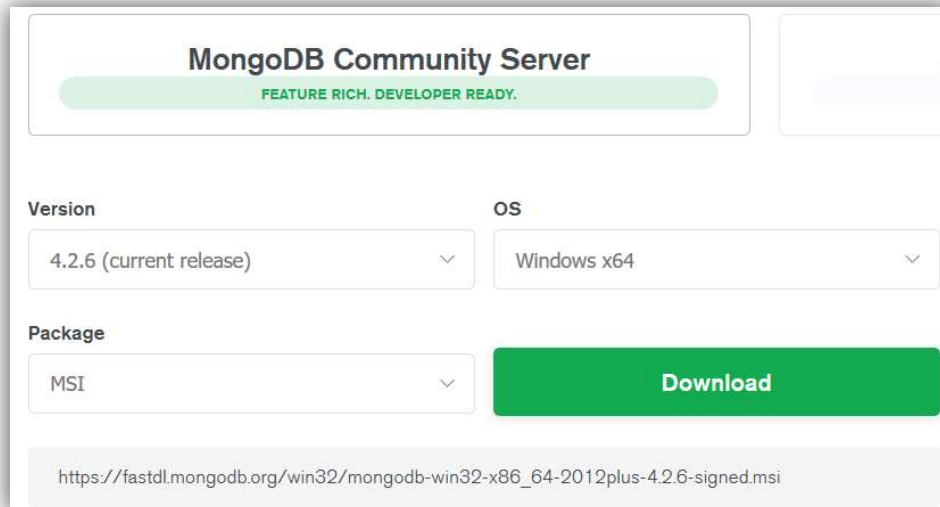
فهرست گزارش سوالات

- بخش اول: ذخیره اطلاعات بدون ساختار / کار با MongoDB 3
- بخش دوم: مدل سازی داده ها با گراف / کار با دیتابیس Neo4j 29
- بخش سوم: دیتابیس های سطر گسترده / کار با HBase 44

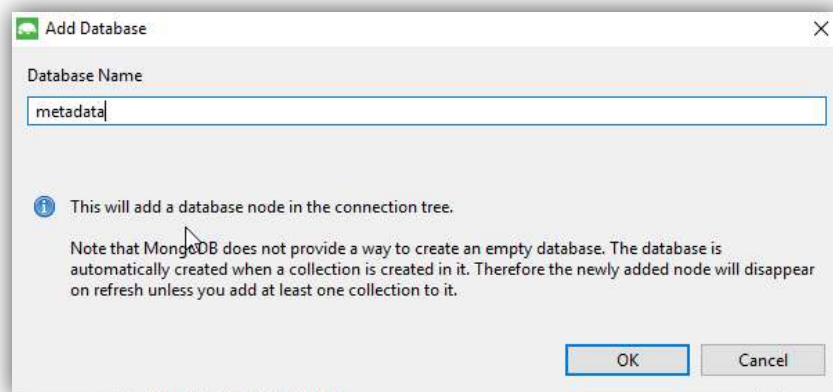
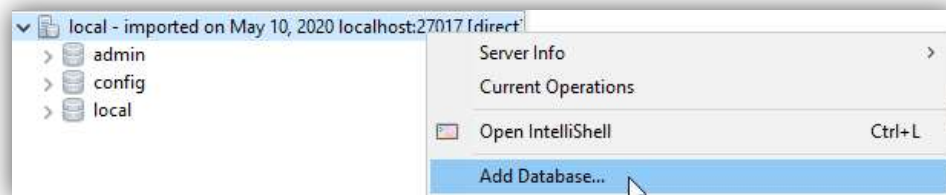
بخش اول: ذخیره اطلاعات بدون ساختار / کار با MongoDB

مراحل دانلود، نصب و آماده‌سازی کار با MongoDB به‌صورت زیر می‌باشد:

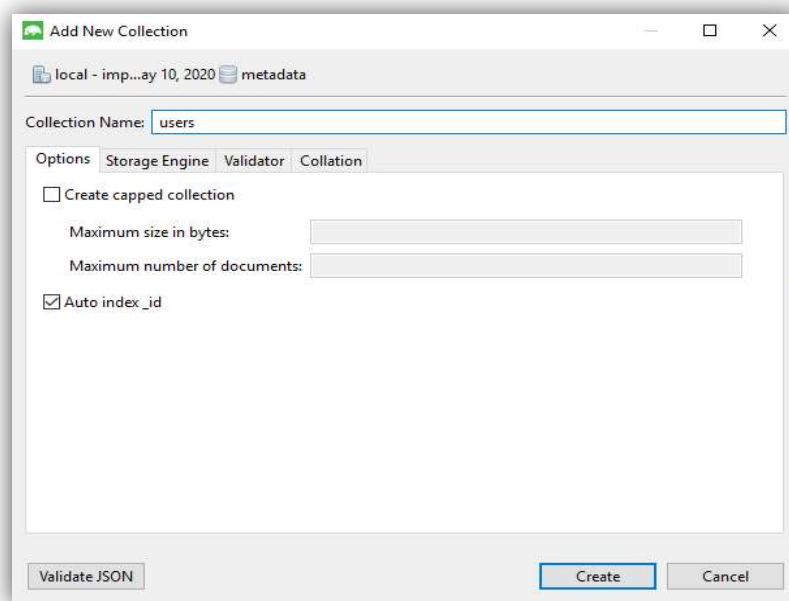
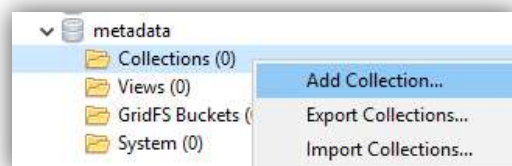
- 1- ابتدا از آدرس <https://www.mongodb.com/download-center/community>، مطابق با تصویر زیر، نسخه‌ی 4.2.6 مانگودی‌بی را به‌صورت بسته‌بندی MSI برای سیستم‌عامل ویندوز دانلود می‌کنیم:



- 2- سپس در آدرس پیش‌فرض خود که همان **C:\Program Files\MongoDB\Server\4.2** است، آنرا نصب می‌کنیم.
- 3- به‌طور پیش‌فرض داده‌های مانگو در آدرس **C:\data\db** ذخیره خواهند شد، به‌همین دلیل، در درایو C یک فولدر به نام data را ایجاد می‌کنیم و سپس در داخل آن یک فولدر به نام db می‌سازیم تا داده‌های مانگو بدون مشکل در این فولدر ذخیره شوند.
- 4- سپس آدرس **C:\Program Files\MongoDB\Server\4.2\bin** را به Path های سیستم اضافه می‌کنیم.
- 5- در نهایت با اجرای دستور mongod.exe در cmd، مانگودی‌بی را اجرا می‌کنیم. در ضمن اگر بخواهیم که از خط فرمان مانگودی‌بی (mongo shell) استفاده کنیم، دستور mongo.exe را در cmd اجرا می‌کنیم.
- 6- در اینجا به‌جای خط فرمان مانگودی‌بی از ابزار گرافیکی بسیار قدرتمند Studio 3T استفاده می‌کنیم که نسخه‌ی قدرتمندتر ابزار Robo 3T می‌باشد.
- 7- در Studio 3T، ابتدا برای اتصال به DB Server به آدرس **localhost:27017** متصل می‌شویم. سپس دیتابیس metadata را به‌صورت زیر ایجاد می‌کنیم:




پس از ایجاد دیتابیس metadata، کالکشن users را به دیتابیس metadata اضافه می‌نماییم:



گام اول تمرین – درج داده‌ها



ابتدا گزینه‌ی  را انتخاب کرده تا به خط فرمان مانگودی‌بی دسترسی داشته باشیم. سپس آدرس <https://randomuser.me/api/?nat=ir> را فراخوانی کرده، فایل JSON آنرا کپی می‌کنیم و سپس اطلاعات موجود در فیلد results را با دستور زیر در کالکشن users وارد می‌کنیم:

```
db.getCollection("users").insertOne({"gender":"female","name":{"title":"Ms","first":  
"الینا","last":"نکو نظر"},"location":{"street":{"number":288,"name":"شهید محمد  
منتظری"},"city":"اهواز","state":"مرکزی","country":"Iran","postcode":32571,"coordinates"  
":{"latitude":45.7234,"longitude":100.7914},"timezone":{"offset":"-  
11:00"},"description":"Midway Island,  
Samoa"}}, {"email":"lyn.nkwnzr@example.com","login":{"uuid":"3bd0072a-4a54-47da-b025-  
a26f3eb9a074","username":"goldenmouse258","password":"boat","salt":"OvrEY09c","md5":  
"0eb8b489793ee4c7dd78dd0dbefb833","sha1":"2313446429c10d1a494ddc8e7b363843fb8de37e"  
,"sha256":"f0119ba3f1166cdbf1e53305cc9131d7d06bd9f9a893ef1ae9e41e3ce4d39861"},"dob":  
{"date":"1967-03-23T19:18:51.750Z","age":53},"registered":{"date":"2019-07-  
16T03:50:04.010Z","age":1},"phone":"086-56042437","cell":"0930-454-  
9056","id":{"name":"","value":null},"picture":{"large":"https://randomuser.me/api/po  
rtraits/women/18.jpg","medium":"https://randomuser.me/api/portraits/med/women/18.jpg",  
,"thumbnail":"https://randomuser.me/api/portraits/thumb/women/18.jpg"},"nat":"IR"})
```

همین عمل را سه بار تکرار کرده و در نهایت از کالکشن users خروجی می‌گیریم (در صفحه‌ی بعد). همانطور که مشاهده می‌شود، به ازای هر insertی که داشته‌ایم، یک ObjectId با نام _id به‌عنوان کلید به مشخصات کاربر اضافه می‌شود.

Key	Value	Type
▼ (1) { _id : 5ec59f1f0a06be7297040af7 }	{ 13 fields }	Document
_id	5ec59f1f0a06be7297040af7	ObjectId
gender	female	String
> name	{ 3 fields }	Object
> location	{ 7 fields }	Object
email	lyn.nkwnzr@example.com	String
> login	{ 7 fields }	Object
> dob	{ 2 fields }	Object
> registered	{ 2 fields }	Object
phone	086-56042437	String
cell	0930-454-9056	String
> id	{ 2 fields }	Object
> picture	{ 3 fields }	Object
nat	IR	String
▼ (2) { _id : 5ec59f490a06be7297040af8 }	{ 13 fields }	Document
_id	5ec59f490a06be7297040af8	ObjectId
gender	female	String
> name	{ 3 fields }	Object
> location	{ 7 fields }	Object
email	ysmyn.hsny@example.com	String
> login	{ 7 fields }	Object
> dob	{ 2 fields }	Object
> registered	{ 2 fields }	Object
phone	019-67880586	String
cell	0989-573-9004	String
> id	{ 2 fields }	Object
> picture	{ 3 fields }	Object
nat	IR	String
▼ (3) { _id : 5ec59f6f0a06be7297040af9 }	{ 13 fields }	Document
_id	5ec59f6f0a06be7297040af9	ObjectId
gender	female	String
> name	{ 3 fields }	Object
> location	{ 7 fields }	Object
email	mhy.khmrw@example.com	String
> login	{ 7 fields }	Object
> dob	{ 2 fields }	Object
> registered	{ 2 fields }	Object
phone	047-82905330	String
cell	0909-153-3197	String
> id	{ 2 fields }	Object
> picture	{ 3 fields }	Object
nat	IR	String

سپس اطلاعات قبل را پاک کرده و با استفاده از کتابخانه‌ی pymongo و به‌صورت زیر، اطلاعات 100000 کاربر تصادفی را در مانگو ذخیره می‌کنیم:

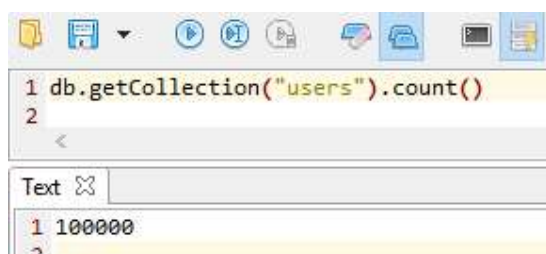
```
from pymongo import MongoClient
import requests, json, time
url = "https://randomuser.me/api/"

client = MongoClient("mongodb://localhost:27017/")
database = client["metadata"]
collection = database["users"]

def get_random_users(n):
    for i in range(0,n//5000):
        querystring = {"nat": "ir", "results": "5000"}
        response = None
        while True:
            response = requests.request("GET", url, params=querystring)
            if response.status_code==503:
                print("120 secs wait...")
                time.sleep(120)
            else:
                break
        data = json.loads(response.text)
        collection.insert_many(data['results'])
        print(i+1)
    print("finished.")

get_random_users(n=100000)
```

در نهایت، با دستور count مطمئن می‌شویم که 100000 کاربر تصادفی در کالکشن users ذخیره شده باشد:



گام دوم تمرین – دستورات اصلی

1. نام و نام خانوادگی کاربرانی که بالای 50 سال سن داشته و ساکن نیشابور باشند:

کوئری:

```
1 db.getCollection("users").find(
2   {
3     "dob.age":{
4       "$gte": 50
5     },
6     "location.city": "نیشابور"
7   },
8   {
9     "name.first": 1,
10    "name.last" : 1,
11    "_id" : 0
12  }
13 )
```

خروجی:

Tree		JSON
Key	Value	
✓ (1) { _id : }	{ 1 fields }	1 {
✓ (1) name	{ 2 fields }	2 "age" : {
"_" first	پارسا	3 "first" : "پارسا",
"_" last	حسینی	4 "last" : "حسینی"
(2) { _id : }	{ 1 fields }	5 }
✓ (1) name	{ 2 fields }	6 }
"_" first	مهرسا	7 {
"_" last	کامیاران	8 "age" : {
(3) { _id : }	{ 1 fields }	9 "first" : "مهرسا",
✓ (1) name	{ 2 fields }	10 "last" : "کامیاران"
"_" first	آرمین	11 }
"_" last	کوئی	12 }
		13 {
		14 "age" : {
		15 "first" : "آرمین",
		16 "last" : "کوئی"
		17 }
		18 }

* لیست کامل خروجی که شامل 995 مورد است، در آدرس [Code\Part 1- MongoDB\Step 2 - main commands](#) قرار گرفته است.

0.191s

زمان اجرا:

2. نام خانوادگی، آدرس و شماره موبایل کاربرانی که بیش از 20 سال است که در سایت ما ثبت نام کرده اند:

نکته: با توجه به اینکه کاربران این سایت، حداکثر 18 سال پیش در این سایت ثبت نام کرده اند، در نتیجه برای اینکه سوال ذکر شده، خروجی داشته باشد، به جای 20، از 18 استفاده خواهیم کرد.

کوئری:

```
1 db.getCollection("users").find(
2   {
3     "registered.age":{
4       "$gte": 18
5     }
6   },
7   {
8     "name.last" : 1,
9     "location": 1,
10    "cell": 1,
11    "_id": 0
12  }
13 )
```

خروجی:

Tree		JSON
Key	Value	
▼ (1) { _id : }	{ 3 fields }	1 {
▼ (1) name	{ 1 fields }	2 "name" : {
"-" last	کریمی	3 "last" : "کریمی"
▼ (1) location	{ 7 fields }	4 },
▼ (1) street	{ 2 fields }	5 "location" : {
number	8352	6 "street" : {
"-" name	پارک لاله	7 "number" : NumberInt(8352),
"-" city	بابل	8 "name" : "پارک لاله"
"-" state	یزد	9 },
"-" country	Iran	10 "city" : "بابل",
postcode	69540	11 "state" : "یزد",
▼ (1) coordinates	{ 2 fields }	12 "country" : "Iran",
latitude	-25.4863	13 "postcode" : NumberInt(69540),
longitude	-116.6241	14 "coordinates" : {
▼ (1) timezone	{ 2 fields }	15 "latitude" : "-25.4863",
offset	-11:00	16 "longitude" : "-116.6241"
description	Midway Island, Samoa	17 },
cell	0938-081-3563	18 "timezone" : {
		19 "offset" : "-11:00",
		20 "description" : "Midway Island, Samoa"
		21 },
		22 "cell" : "0938-081-3563"
		23 }
		24 }

* لیست کامل خروجی که شامل 4399 مورد است، در آدرس [Code\Part 1- MongoDB\Step 2 - main commands](#) قرار گرفته است.

0.219s

زمان اجرا:

3. افزودن فیلد year_persian به فیلد dob و registered:

کوئری:

```
1 db.getCollection("users").update(
2   {
3     },
4   [
5     {
6       "$set": {
7         "dob.year_persian": {
8           "$year": {
9             "$subtract": [ {"$dateFromString": {"dateString": "$dob.date"}}, 226899*24*60*60000 ]
10          }
11        },
12        "registered.year_persian": {
13          "$year": {
14            "$subtract": [ {"$dateFromString": {"dateString": "$registered.date"}}, 226899*24*60*60000 ]
15          }
16        }
17      }
18    },
19    {
20      "multi": true
21    }
22  ]
23 )
```

توضیح: ابتدا نوع فیلد dob.date که به صورت رشته ذخیره شده است را به Date تغییر می دهیم. سپس با توجه به اینکه اختلاف تاریخ شمسی و میلادی 226899 روز می باشد، 226899 روز از این تاریخ کم می کنیم تا به تاریخ شمسی دست یابیم. در نهایت با استفاده از \$year، سال را از تاریخ شمسی حاصل استخراج می کنیم و آنرا به فیلد جدیدی به نام dob.year_persian اختصاص می دهیم. ما این عمل را برای registered.date نیز انجام می دهیم.

خروجی کالکشن users به صورت Tree:

Key	Value
<ul style="list-style-type: none"> <ul style="list-style-type: none"> { 13 fields } _id: 5ec5aec0ab89eef4d3aaa3c2 gender: female name: { 3 fields } location: { 7 fields } email: lyn.hmdy@example.com login: { 7 fields } dob: { 3 fields } <ul style="list-style-type: none"> date: 1988-07-07T06:26:12.914Z age: 32 year_persian: 1367 registered: { 3 fields } <ul style="list-style-type: none"> date: 2005-12-13T16:08:51.103Z age: 15 year_persian: 1384 phone: 031-23157485 cell: 0945-033-0974 id: { 2 fields } picture: { 3 fields } nat: IR 	

* به همین صورت تمام ردیف های کالکشن users ویرایش شده، حاوی زیرفیلد year_persian در فیلد dob و registered هستند. با توجه به اینکه خروجی این سوال، حدود 180 مگابایت حجم دارد، از ارائه ی آن صرف نظر می کنیم.

12.03s

زمان اجرا:

4. دستوری بنویسید که با اجرا شدن آن در هر روز، نام و نام خانوادگی و ایمیل افرادی که در آن روز تولدشان است، به ما برگردانده شود:

کوئری:

```
1 db.getCollection("users").find(
2   {
3     "$expr":{
4       "$and":[
5         {
6           "$eq":[
7             {"$month":{"date":{"$dateFromString":{"dateString":"$dob.date"}}}},
8             {"$month":{"date:new Date()}}
9           ]
10        },
11        {
12          "$eq":[
13            {"$dayOfMonth":{"date":{"$dateFromString":{"dateString":"$dob.date"}}}},
14            {"$dayOfMonth":{"date:new Date()}}
15          ]
16        }
17      ]
18    }
19  },
20  {
21    "name.first": 1,
22    "name.last": 1,
23    "email":1,
24    "_id":0
25  }
26 )
```

توضیح: در کوئری بالا، نام و نام خانوادگی افرادی که روز و ماه تولدشان با روز و ماه جاری برابر است، برگردانده می شود. لازم به ذکر است که میتوانستیم بجای \$month و \$dayOfMonth از \$dayOfYear استفاده کنیم، ولی چونکه با این روش، سال کبیسه در نظر گرفته نمی شود، به همین دلیل از روش اول استفاده نمودیم.

خروجی:

Tree		JSON
Key	Value	
▼ (1) { _id : }	{ 2 fields }	1 {
▼ { name }	{ 2 fields }	2 "name" : {
{ "first" }	مهدی	3 "first" : "مهدی",
{ "last" }	نجاتی	4 "last" : "نجاتی"
{ "email" }	mhdy.njty@example.com	5 },
{ "email" }	mhdy.njty@example.com	6 "email" : "mhdy.njty@example.com"
▼ (2) { _id : }	{ 2 fields }	7 }
▼ { name }	{ 2 fields }	8 {
{ "first" }	محمدپارسا	9 "name" : {
{ "last" }	سلطانی نژاد	10 "first" : "محمدپارسا",
{ "email" }	mhmdprs.sltnynjd@example.com	11 "last" : "سلطانی نژاد"
		12 },
		13 "email" : "mhmdprs.sltnynjd@example.com"
		14 }

* لیست کامل خروجی که شامل 251 مورد است، در آدرس [Code\Part 1- MongoDB\Step 2 - main commands](#) قرار گرفته است.

3.041s

زمان اجرا:

5. ذخیره پسورد به شکل خام در اطلاعات کاربر، کاری غیر حرفه‌ای است. می‌خواهیم این مشکل را برطرف کنیم. چه راه حلی برای حل مساله پیشنهاد میکنید؟
برای حل این مسئله، می‌توان بجای پسورد از hex_md5 آن استفاده کرد.

راه حل را ابتدا روی یک سند خاص، امتحان کنید و مطمئن شوید بعد از اعمال تغییرات، آن کاربر خاص را با داشتن یوزنیم و پسورد، میتواند پیدا کند. سپس تمام کاربران را به روز رسانی کنید.

باتوجه به اینکه hex_md5 یک تابع جاوا اسکریپتی است، در نتیجه می‌توان از این تابع در خط فرمان مانگودی بی استفاده کرد. ولی بایستی ورودی این تابع، پسورد کاربر مورد نظر باشد. ولی مشکلی که وجود دارد این است که در یک تابع جاوا اسکریپتی، نمی‌توان با "\$login.password"، مقدار پسورد کاربر مورد نظر را دریافت و از آن به عنوان ورودی تابع hex_md5 استفاده کرد. به همین دلیل ما مجبور شدیم که از foreach بر روی find استفاده کنیم تا بتوانیم با دستورات جاوا اسکریپتی، پسورد کاربر(های) مورد نیاز را استخراج و از آنها به عنوان ورودی تابع hex_md5 استفاده کنیم. لازم به ذکر است که با توجه به اینکه استفاده از دستورات و توابع جاوا اسکریپتی در مانگودی بی مجاز است، در نتیجه تمام دستورات زیر در خط فرمان مانگودی بی قابل اجرا است.

ابتدا با استفاده از کوئری زیر، به جای پسورد کاربری با نام کاربری smallrabbit981، مقدار hex_md5 آنرا قرار می‌دهیم:

```
db.getCollection("users").find({"login.username":"smallrabbit981"}).forEach(function(data){
    data.login.password = hex_md5(data.login.password)
    db.getCollection("users").save(data)
});
```

پسورد خام کاربر smallrabbit981 ، daytona بود. حال اگر بخواهیم که این کاربر را با داشتن نام کاربری و پسورد خام پیدا کنیم، کفایت از کوئری زیر استفاده کنیم:

```
1 db.getCollection("users").find(
2   {
3     "login.username":"smallrabbit981",
4     "login.password":hex_md5("daytona")
5   }
6 )
```

که خروجی آن به صورت زیر می‌باشد:

Key	Value
▼ (1) { _id : 5ec5aec0ab89eef4d3aaa3a4 }	{ 13 fields }
_id	5ec5aec0ab89eef4d3aaa3a4
gender	male
name	{ 3 fields }
location	{ 7 fields }
email	prs.sly@example.com
login	{ 7 fields }
uuid	e4e3ff7f-d6e5-487b-84dc-f8e3c921eebb
username	smallrabbit981
password	5d1fb338b351a63a1720b837cc9010ff
salt	2x5wGESI
md5	7057c28225012a5df1d82ba343b5dcc0
sha1	ea9ee870538e7e7bf5ea85af770d2713937672
sha256	6a4ec63fba4177a656cc53d9e5683009f13367
dob	{ 3 fields }
registered	{ 3 fields }
phone	075-11897642
cell	0919-951-5198
id	{ 2 fields }
picture	{ 3 fields }
nat	IR

با توجه به اینکه مطمئن شدیم که عملکرد روش ما درست است، در نتیجه با کوئری زیر، پسورد تمام کاربران را به روز رسانی می‌کنیم (البته قبل از آن، پسورد کاربر بالا را به صورت قبل برمی‌گردانیم):

```
db.getCollection("users").find().forEach(function(data){
  data.login.password = hex_md5(data.login.password)
  db.getCollection("users").save(data)
});
```

همانطور که از خروجی username و password چند کاربر اول کالکشن users به صورت زیر مشخص است، پسورد کاربران تغییر کرده است:

```
1 {
2   "login" : {
3     "username" : "smallrabbit981",
4     "password" : "5d1fb338b351a63a1720b837cc9010ff"
5   }
6 }
7 {
8   "login" : {
9     "username" : "beautifulmeercat896",
10    "password" : "507fcbb50bf9e88168dcd0d083600c0b"
11  }
12 }
13 {
14   "login" : {
15     "username" : "whiterabbit388",
16     "password" : "7bf4a177e86772277c63b6d5352741"
17   }
18 }
```

زمان اجرای کوئری به روز رسانی پسورد تمام کاربران نیز 2 دقیقه و 2 ثانیه و 50 صدم ثانیه به طول انجامید.

گام سوم تمرین – دستورات تجمعی و آماری

1. ابتدا بایستی کاربران را به سه گروه سنی کمتر از 16 (نوجوان) / بین 16 تا 30 (جوان) / بالاتر از 30 (میانسال به بالا) تقسیم کنید، سپس دستوری بنویسید که تعداد هر گروه را به ما برگرداند:

کوئری:

```
1 db.getCollection("users").aggregate(
2   [
3     {
4       "$project" : {
5         "AgeGroup" : {
6           "$cond" : {
7             "if" : { "$lt" : [ "$dob.age", 16.0 ] },
8             "then" : "teenager",
9             "else" : {
10              "$cond" : {
11                "if" : { "$lte" : [ "$dob.age", 30.0 ] },
12                "then" : "young",
13                "else" : "middle-aged"
14              }
15            }
16          }
17        }
18      },
19    ],
20    {
21      "$group" : {
22        "_id" : "$AgeGroup",
23        "count" : { "$sum" : 1.0 }
24      }
25    },
26    {
27      "$project" : {
28        "_id" : 0.0,
29        "age-group" : "$_id",
30        "count" : { "$toInt" : "$count" }
31      }
32    }
33  ]
34 )
```

خروجی:

Tree		JSON
Key	Value	
▼ (1) { _id : }	{ 2 fields }	1 {
age-group	young	2 "age-group" : "young",
count	16006	3 "count" : NumberInt(16006)
		4 }
▼ (2) { _id : }	{ 2 fields }	5 {
age-group	middle-aged	6 "age-group" : "middle-aged",
count	83994	7 "count" : NumberInt(83994)
		8 }

0.389s

زمان اجرا:







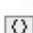



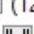


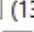


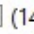


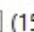

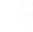



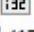
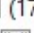


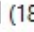





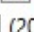

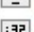
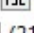


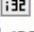


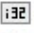


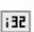
2. تعداد کاربران هر استان را به تفکیک تولید کنید:

کوئری:

```
1 db.getCollection("users").aggregate(
2   [
3     {
4       "$group" : {
5         "_id" : "$location.state",
6         "count" : { "$sum" : 1.0 }
7       }
8     },
9     {
10      "$project" : {
11        "_id" : 0.0,
12        "state" : "$_id",
13        "count" : { "$toInt" : "$count" }
14      }
15    }
16  ]
17 )
```

خروجی:

Tree		JSON
Key	Value	
▼ (1) { _id : }	{ 2 fields }	1 {
state	قم	2 "state" : "قم",
count	3253	3 "count" : NumberInt(3253)
▼ (2) { _id : }	{ 2 fields }	4 }
state	آذربایجان غربی	5 {
count	3186	6 "state" : "آذربایجان غربی",
▼ (3) { _id : }	{ 2 fields }	7 "count" : NumberInt(3186)
state	سیستان و بلوچستان	8 }
count	3239	9 {
▼ (4) { _id : }	{ 2 fields }	10 "state" : "سیستان و بلوچستان",
state	خراسان رضوی	11 "count" : NumberInt(3239)
count	3324	12 }
▼ (5) { _id : }	{ 2 fields }	13 {
state	اردبیل	14 "state" : "خراسان رضوی",
count	3187	15 "count" : NumberInt(3324)
▼ (6) { _id : }	{ 2 fields }	16 }
state	خراسان شمالی	17 {
count	3196	18 "state" : "اردبیل",
▼ (7) { _id : }	{ 2 fields }	19 "count" : NumberInt(3187)
state	لرستان	20 }
count	3205	21 {
▼ (8) { _id : }	{ 2 fields }	22 "state" : "خراسان شمالی",
state	اصفهان	23 "count" : NumberInt(3196)
count	3244	24 }
		25 {
		26 "state" : "لرستان",
		27 "count" : NumberInt(3205)
		28 }
		29 {
		30 "state" : "اصفهان",
		31 "count" : NumberInt(3244)
		32 }

▼  (9) { _id : }	{ 2 fields }	33 {	
 state	کهگیلویه و بویراحمد	34 "state" : "کهگیلویه و بویراحمد",	
 count	3116	35 "count" : NumberInt(3116)	
▼  (10) { _id : }	{ 2 fields }	36 }	
 state	یزد	37 {	
 count	3262	38 "state" : "یزد",	
▼  (11) { _id : }	{ 2 fields }	39 "count" : NumberInt(3262)	
 state	گلستان	40 }	
 count	3297	41 {	
▼  (12) { _id : }	{ 2 fields }	42 "state" : "گلستان",	
 state	کرمانشاه	43 "count" : NumberInt(3297)	
 count	3178	44 }	
▼  (13) { _id : }	{ 2 fields }	45 {	
 state	هرمزگان	46 "state" : "کرمانشاه",	
 count	3261	47 "count" : NumberInt(3178)	
▼  (14) { _id : }	{ 2 fields }	48 }	
 state	ایلام	49 {	
 count	3356	50 "state" : "هرمزگان",	
▼  (15) { _id : }	{ 2 fields }	51 "count" : NumberInt(3261)	
 state	مازندران	52 }	
 count	3290	53 {	
▼  (16) { _id : }	{ 2 fields }	54 "state" : "ایلام",	
 state	چهارمحال و بختیاری	55 "count" : NumberInt(3356)	
 count	3206	56 }	
▼  (17) { _id : }	{ 2 fields }	57 {	
 state	خراسان جنوبی	58 "state" : "مازندران",	
 count	3205	59 "count" : NumberInt(3290)	
▼  (18) { _id : }	{ 2 fields }	60 }	
 state	گیلان	61 {	
 count	3268	62 "state" : "چهارمحال و بختیاری",	
▼  (19) { _id : }	{ 2 fields }	63 "count" : NumberInt(3206)	
 state	مرکزی	64 }	
 count	3205	65 {	
▼  (20) { _id : }	{ 2 fields }	66 "state" : "خراسان جنوبی",	
 state	همدان	67 "count" : NumberInt(3205)	
 count	3207	68 }	
▼  (21) { _id : }	{ 2 fields }	69 {	
 state	آذربایجان شرقی	70 "state" : "گیلان",	
 count	3185	71 "count" : NumberInt(3268)	
▼  (22) { _id : }	{ 2 fields }	72 }	
 state	بوشهر	73 {	
 count	3170	74 "state" : "مرکزی",	
▼  (23) { _id : }	{ 2 fields }	75 "count" : NumberInt(3205)	
 state	کردستان	76 }	
 count	3241	77 {	
▼  (24) { _id : }	{ 2 fields }	78 "state" : "همدان",	
 state	خوزستان	79 "count" : NumberInt(3207)	
 count	3222	80 }	
		81 {	
		82 "state" : "آذربایجان شرقی",	
		83 "count" : NumberInt(3185)	
		84 }	
		85 {	
		86 "state" : "بوشهر",	
		87 "count" : NumberInt(3170)	
		88 }	

0.247s

17

3. فیلد شماره موبایل افرادی که آفست timezone آنها برابر +5:00 را حذف نمایید، سپس تعداد افرادی را بیابید که شماره موبایل ندارند.

کوئری حذف شماره موبایل:

```
1 db.getCollection("users").update(
2   {
3     "location.timezone.offset": "+5:00"
4   },
5   [{
6     $unset: "cell"
7   }],
8   {
9     "multi": true
10  }
11 )
```

کوئری جست و جوی تعداد افراد:

```
1 db.getCollection("users").aggregate(
2   [
3     {
4       "$match" : {
5         "cell" : { "$exists" : false }
6       }
7     },
8     {
9       "$group" : {
10        "_id" : null,
11        "count" : { "$sum" : 1 }
12      }
13    },
14    {
15      "$project" : {
16        "_id" : 0,
17        "count" : { "$toInt" : "$count" }
18      }
19    }
20  ]
21 )
```

خروجی:

Tree		JSON
Key	Value	1 {
▼ (1) { _id : }	{ 1 fields }	2 "count" : NumberInt(3398)
count	3398	3 }

0.096s

زمان اجرا:

4. میانگین سن کاربران استان تهران را با میانگین سن کاربران سایر استان ها مقایسه کنید:

کوئری:

درک اینجانب از سوال، این است که میانگین سن کاربران استان تهران را با میانگین سن دیگر کاربران مقایسه کنیم. به همین دلیل، کوئری ما به صورت زیر می باشد:

```
1 db.getCollection("users").aggregate(  
2   [  
3     {  
4       "$project" : {  
5         "state" : {  
6           "$cond" : {  
7             "if" : { "$eq" : [ "$location.state", "تهران" ] },  
8             "then" : "تهران",  
9             "else" : "سایر"  
10          }  
11        },  
12        "age" : "$dob.age"  
13      }  
14    },  
15    {  
16      "$group" : {  
17        "_id" : "$state",  
18        "avgAge" : { "$avg" : "$age"}  
19      }  
20    },  
21    {  
22      "$project" : {  
23        "_id" : 0.0,  
24        "state" : "$_id",  
25        "avgAge" : "$avgAge"  
26      }  
27    }  
28  ]  
29 )
```

خروجی:

Tree		JSON
Key	Value	
✓ (1) { _id : }	{ 2 fields }	1 {
"state"	سایر	2 "state" : "سایر",
avgAge	48.87177288552331	3 "avgAge" : 48.87177288552331
✓ (2) { _id : }	{ 2 fields }	4 }
"state"	تهران	5 {
avgAge	48.39993755853887	6 "state" : "تهران",
		7 "avgAge" : 48.39993755853887
		8 }

همانطور که مشاهده می شود، میانگین سن کاربران استان تهران از میانگین سن دیگر کاربران کمتر است.

0.425s

زمان اجرا:

کوئری: ولی اگر هدف سوال از مقایسه این است که مشخص کنیم که چه استان‌هایی، میانگین سن کمتر، چه استان‌هایی، میانگین سن برابر و چه استان‌هایی، میانگین سن بیشتر از استان تهران را دارند، کوئری آن به صورت زیر می‌باشد:

```

1 db.getCollection("users").aggregate(
2   [
3     {
4       "$group" : {
5         "_id" : "$location.state",
6         "avgAge" : { "$avg" : "$dob.age" }
7       }
8     },
9     {
10      "$facet" : {
11        "tehran" : [
12          {
13            "$match" : { "$expr" : { "$eq" : [ "$_id", "تهران" ] }}
14          },
15          {
16            "$project" : { "avg" : "$avgAge" }
17          }
18        ],
19        "states" : [
20          {
21            "$match" : { "$expr" : { "$ne" : [ "$_id", "تهران" ] }}
22          },
23          {
24            "$project" : { "_id" : 1.0, "avgAge" : 1.0 }
25          }
26        ]
27      },
28    },
29    {
30      "$facet" : {
31        "less" : [
32          {
33            "$match" : { "$expr" : { "$lt" : [ "$states.avgAge", "$tehran.avg" ] }}
34          },
35          {
36            "$group" : { "_id" : "$states", "states" : { "$addToSet" : "$states._id" }}
37          }
38        ],
39        "eq1" : [
40          {
41            "$match" : { "$expr" : { "$eq" : [ "$states.avgAge", "$tehran.avg" ] }}
42          },
43          {
44            "$group" : { "_id" : "$states", "states" : { "$addToSet" : "$states._id" }}
45          }
46        ],
47        "grt" : [
48          {
49            "$match" : { "$expr" : { "$gt" : [ "$states.avgAge", "$tehran.avg" ] }}
50          },
51          {
52            "$group" : { "_id" : "$states", "states" : { "$addToSet" : "$states._id" }}
53          }
54        ]
55      },
56    },
57    {
58      "$project" : {
59        "less_cities" : { "$arrayElemAt" : [{ "$arrayElemAt" : [ "$less.states", 0.0 ], 0.0 }],
60        "equal_cities" : { "$arrayElemAt" : [{ "$arrayElemAt" : [ "$eq1.states", 0.0 ], 0.0 }],
61        "greater_cities" : { "$arrayElemAt" : [{ "$arrayElemAt" : [ "$grt.states", 0.0 ], 0.0 }],
62      }
63    }
64  ]
65 )

```

خروجی:

Tree		JSON
Key	Value	
▼ { } (1) { _id : }	{ 3 fields }	1 {
less_cities	null	2 "less_cities" : null,
equal_cities	null	3 "equal_cities" : null,
▼ [] greater_cities	[30 elements]	4 "greater_cities" : [
0	سمنان	5 "سمنان",
1	البرز	6 "البرز",
2	قزوین	7 "قزوین",
3	فارس	8 "فارس",
4	زنجان	9 "زنجان",
5	کرمان	10 "کرمان",
6	خوزستان	11 "خوزستان",
7	بوشهر	12 "بوشهر",
8	آذربایجان شرقی	13 "آذربایجان شرقی",
9	همدان	14 "همدان",
10	کردستان	15 "کردستان",
11	گیلان	16 "گیلان",
12	مرکزی	17 "مرکزی",
13	خراسان جنوبی	18 "خراسان جنوبی",
14	چهارمحال و بختیاری	19 "چهارمحال و بختیاری",
15	مازندران	20 "مازندران",
16	هرمزگان	21 "هرمزگان",
17	ایلام	22 "ایلام",
18	کرمانشاه	23 "کرمانشاه",
19	گلستان	24 "گلستان",
20	یزد	25 "یزد",
21	کهگیلویه و بویراحمد	26 "کهگیلویه و بویراحمد",
22	اصفهان	27 "اصفهان",
23	لرستان	28 "لرستان",
24	خراسان شمالی	29 "خراسان شمالی",
25	خراسان رضوی	30 "خراسان رضوی",
26	اردبیل	31 "اردبیل",
27	سیستان و بلوچستان	32 "سیستان و بلوچستان",
28	آذربایجان غربی	33 "آذربایجان غربی",
29	قم	34 "قم",
		35]
		36 }

همانطور که مشاهده می‌شود، میانگین سن کاربران استان تهران از میانگین سن کاربران تمام استان‌های دیگر کمتر است.

0.317s

زمان اجرا:

کوئری:

در اینجا، ابتدا تعداد کاربران هر استان را محاسبه کرده، سپس یک بار، استان‌ها را بر اساس تعداد کاربرانشان به‌صورت نزولی مرتب کرده و اولین استان را برمی‌داریم، و یک بار، استان‌ها را بر اساس تعداد کاربرانشان به‌صورت صعودی مرتب کرده و اولین استان را برمی‌داریم. سپس این دو استان را به‌صورت مناسب در خروجی نمایش می‌دهیم.

خروجی:

0.250s


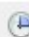
زمان اجرا:

گام چهارم تمرین - بررسی کارآیی شاخص‌ها


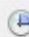
به منظور ایجاد شاخص، مراحل زیر در ابزار Studio 3T طی می‌شود:
ابتدا بر روی کالکشن users کلیک راست نموده، سپس گزینه‌ی Add Index را انتخاب می‌کنیم.
حال از طریق گزینه‌ی Add Field(s)، فیلد(های) مورد نظر را انتخاب می‌نماییم.
پس از انتخاب فیلد(های) مورد نظر، Index Type آنها را از بین یکی از موارد زیر انتخاب می‌کنیم:

ascending
descending
hashed
text
2dsphere
2d
geoHaystack

استفاده از شاخص در گام دوم - سوال 1 :

اجرای کوئری با شاخص	اجرای کوئری بدون شاخص	Indexed Fields:	
 0.014s	 0.191s	Field	Index Type
		dob.age	descending
		location.city	descending

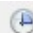
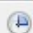
استفاده از شاخص در گام دوم - سوال 4:

اجرای کوئری با شاخص	اجرای کوئری بدون شاخص	Indexed Fields: <table><tr><th>Field</th><th>Index Type</th></tr><tr><td>dob.date</td><td>ascending</td></tr></table>	Field	Index Type	dob.date	ascending
Field	Index Type					
dob.date	ascending					
 2.853s	 3.041s					



استفاده از شاخص در گام دوم - سوال 5:

اجرای کوئری با شاخص	اجرای کوئری بدون شاخص	Indexed Fields:	
1 دقیقه و 59 ثانیه و 8 صدم ثانیه	2 دقیقه و 2 ثانیه و 50 صدم ثانیه	Field	Index Type
		login.password	ascending

استفاده از شاخص در گام سوم - سوال 2 :

اجرای کوئری با شاخص	اجرای کوئری بدون شاخص	<div>Indexed Fields:</div> <table><tr><th>Field</th><th>Index Type</th></tr><tr><td>location.state</td><td>ascending</td></tr></table>	Field	Index Type	location.state	ascending
Field	Index Type					
location.state	ascending					
<div> 0.242s</div>	<div> 0.247s</div>					

استفاده از شاخص در گام سوم – سوال 5:

اجرای کوئری با شاخص	اجرای کوئری بدون شاخص	<div>Indexed Fields:</div> <table><tr><th>Field</th><th>Index Type</th></tr><tr><td>location.city</td><td>ascending</td></tr></table>	Field	Index Type	location.city	ascending
Field	Index Type					
location.city	ascending					
<div> 0.272s</div>	<div> 0.250s</div>					

همانطور که از جداول بالا مشخص شد، استفاده از شاخص در مسئله‌ی گام دوم – سوال 1 در کاهش زمان اجرا بسیار تاثیرگذار بوده است (که دلیل آن این است که فقط از find استفاده شده است). ولی استفاده از شاخص در سایر موارد، تاثیر چندانی در کاهش زمان اجرا نداشته است، که احتمالا دلیل آن به نوع کوئری نوشته شده برمی‌گردد.

گام پنجم تمرین – Map Reduce

اسکرپت زیر، قابلیت بازیابی میانگین سن کاربران به تفکیک هر شهر را دارد. این اسکرپت با توجه به اینکه به زبان جاوا اسکرپت نوشته شده است، در نتیجه می‌تواند در خط فرمان مانگودی‌بی نیز اجرا شود:

```
// *** 3T Software Labs, Studio 3T: MapReduce Job ***

// Variable for db
var __3tsoftwarelabs_db = "metadata";


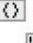

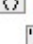
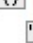
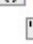



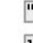
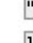
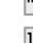
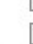


// Variable for map
var __3tsoftwarelabs_map = function () {
    emit(this.location.city, this.dob.age);
};

// Variable for reduce
var __3tsoftwarelabs_reduce = function (key, values) {
    return (key, Array.avg(values));
};

db.runCommand({
    mapReduce: "users",
    map: __3tsoftwarelabs_map,
    reduce: __3tsoftwarelabs_reduce,
    out: {"inline": 1},
    query: {},
    sort: {},
    inputDB: "metadata",
});
```

توضیح: خروجی تابع map، لیستی از زوج مرتب‌های (مقدار، کلید) به صورت (سن، شهر) است. سپس تابع reduce، به ازای هر شهر، میانگین سن کاربران آنها را محاسبه و برمی‌گرداند. با اجرای این اسکرپت در خط فرمان مانگودی‌بی، خروجی آن به صورت زیر خواهد بود:

Tree		JSON
Key	Value	
✓ (1) {_id: آبادان}	{ 2 fields }	1 {
"_id"	آبادان	2 "_id" : "آبادان",
1.23 value	46.30684734774464	3 "value" : 46.30684734774464
4 }		4 }
✓ (2) {_id: آمل}	{ 2 fields }	5 {
"_id"	آمل	6 "_id" : "آمل",
1.23 value	62.7272287416071	7 "value" : 62.7272287416071
8 }		8 }
✓ (3) {_id: اراک}	{ 2 fields }	9 {
"_id"	اراک	10 "_id" : "اراک",
1.23 value	53.29119558235929	11 "value" : 53.29119558235929
12 }		12 }
✓ (4) {_id: اردبیل}	{ 2 fields }	13 {
"_id"	اردبیل	14 "_id" : "اردبیل",
1.23 value	43.21822560199644	15 "value" : 43.21822560199644
16 }		16 }
✓ (5) {_id: ارومیه}	{ 2 fields }	
"_id"	ارومیه	
1.23 value	47.02806895439129	

▼  (6) {_id : اسلام شهر}	{ 2 fields }
"_"_id	اسلام شهر
1.23 value	46.04880682299492
▼  (7) {_id : اصفهان}	{ 2 fields }
"_"_id	اصفهان
1.23 value	40.62441790895408
▼  (8) {_id : اهواز}	{ 2 fields }
"_"_id	اهواز
1.23 value	44.31602943268522
▼  (9) {_id : ایلام}	{ 2 fields }
"_"_id	ایلام
1.23 value	43.454963331113426
▼  (10) {_id : بابل}	{ 2 fields }
"_"_id	بابل
1.23 value	48.43959303896065
▼  (11) {_id : بجنورد}	{ 2 fields }
"_"_id	بجنورد
1.23 value	39.1487795341395
▼  (12) {_id : بروجرد}	{ 2 fields }
"_"_id	بروجرد
1.23 value	38.23225981929865
▼  (13) {_id : بندرعباس}	{ 2 fields }
"_"_id	بندرعباس
1.23 value	48.83462002017418
▼  (14) {_id : بوشهر}	{ 2 fields }
"_"_id	بوشهر
1.23 value	45.187953695153055
▼  (15) {_id : بیرجند}	{ 2 fields }
"_"_id	بیرجند
1.23 value	51.973124742873466
▼  (16) {_id : تبریز}	{ 2 fields }
"_"_id	تبریز
1.23 value	61.177931171109336
▼  (17) {_id : تهران}	{ 2 fields }
"_"_id	تهران
1.23 value	49.52043100284166
▼  (18) {_id : خرم آباد}	{ 2 fields }
"_"_id	خرم آباد
1.23 value	48.11363125019587
▼  (19) {_id : خمینی شهر}	{ 2 fields }
"_"_id	خمینی شهر
1.23 value	48.07379103960122
▼  (20) {_id : خوی}	{ 2 fields }
"_"_id	خوی
1.23 value	38.81108519607666

```

17 {
18   "_id" : "ارومیه",
19   "value" : 47.02806895439129
20 }
21 {
22   "_id" : "اسلام شهر",
23   "value" : 46.04880682299492
24 }
25 {
26   "_id" : "اصفهان",
27   "value" : 40.62441790895408
28 }
29 {
30   "_id" : "اهواز",
31   "value" : 44.31602943268522
32 }
33 {
34   "_id" : "ایلام",
35   "value" : 43.454963331113426
36 }
37 {
38   "_id" : "بابل",
39   "value" : 48.43959303896065
40 }
41 {
42   "_id" : "بجنورد",
43   "value" : 39.1487795341395
44 }
45 {
46   "_id" : "بروجرد",
47   "value" : 38.23225981929865
48 }
49 {
50   "_id" : "بندرعباس",
51   "value" : 48.83462002017418
52 }
53 {
54   "_id" : "بوشهر",
55   "value" : 45.187953695153055
56 }
57 {
58   "_id" : "بیرجند",
59   "value" : 51.973124742873466
60 }
61 {
62   "_id" : "تبریز",
63   "value" : 61.177931171109336
64 }
65 {
66   "_id" : "تهران",
67   "value" : 49.52043100284166
68 }
69 {
70   "_id" : "خرم آباد",
71   "value" : 48.11363125019587
72 }
73 {
74   "_id" : "خمینی شهر",
75   "value" : 48.07379103960122
76 }
77 {
78   "_id" : "خوی",
79   "value" : 38.81108519607666
80 }

```

▼  (21) {_id : دزفول}	{ 2 fields }
"_"_id	دزفول
1.23 value	48.29968213399863
▼  (22) {_id : رشت}	{ 2 fields }
"_"_id	رشت
1.23 value	52.066882332382036
▼  (23) {_id : زاهدان}	{ 2 fields }
"_"_id	زاهدان
1.23 value	37.86827775072374
▼  (24) {_id : زنجان}	{ 2 fields }
"_"_id	زنجان
1.23 value	54.3088577806183
▼  (25) {_id : ساری}	{ 2 fields }
"_"_id	ساری
1.23 value	47.71708878726032
▼  (26) {_id : ساوه}	{ 2 fields }
"_"_id	ساوه
1.23 value	52.74233929800202
▼  (27) {_id : سبزوار}	{ 2 fields }
"_"_id	سبزوار
1.23 value	49.712894559349564
▼  (28) {_id : سنندج}	{ 2 fields }
"_"_id	سنندج
1.23 value	45.54814212605232
▼  (29) {_id : سیرجان}	{ 2 fields }
"_"_id	سیرجان
1.23 value	54.96772413705547
▼  (30) {_id : شهریار}	{ 2 fields }
"_"_id	شهریار
1.23 value	46.573488509275094
▼  (31) {_id : شیراز}	{ 2 fields }
"_"_id	شیراز
1.23 value	47.20046191883396
▼  (32) {_id : قائمشهر}	{ 2 fields }
"_"_id	قائم‌شهر
1.23 value	52.390764656627965
▼  (33) {_id : قدس}	{ 2 fields }
"_"_id	قدس
1.23 value	55.629090385138944
▼  (34) {_id : قرچک}	{ 2 fields }
"_"_id	قرچک
1.23 value	54.74481202576004
▼  (35) {_id : قزوین}	{ 2 fields }
"_"_id	قزوین
1.23 value	41.29597490186893
▼  (36) {_id : قم}	{ 2 fields }
"_"_id	قم
1.23 value	43.728207532391416

```

81 {
82   "_id" : "دزفول",
83   "value" : 48.29968213399863
84 }
85 {
86   "_id" : "رشت",
87   "value" : 52.066882332382036
88 }
89 {
90   "_id" : "زاهدان",
91   "value" : 37.86827775072374
92 }
93 {
94   "_id" : "زنجان",
95   "value" : 54.3088577806183
96 }
97 {
98   "_id" : "ساری",
99   "value" : 47.71708878726032
100 }
101 {
102   "_id" : "ساوه",
103   "value" : 52.74233929800202
104 }
105 {
106   "_id" : "سبزوار",
107   "value" : 49.712894559349564
108 }
109 {
110   "_id" : "سنندج",
111   "value" : 45.54814212605232
112 }
113 {
114   "_id" : "سیرجان",
115   "value" : 54.96772413705547
116 }
117 {
118   "_id" : "شهریار",
119   "value" : 46.573488509275094
120 }
121 {
122   "_id" : "شیراز",
123   "value" : 47.20046191883396
124 }
125 {
126   "_id" : "قائم‌شهر",
127   "value" : 52.390764656627965
128 }
129 {
130   "_id" : "قدس",
131   "value" : 55.629090385138944
132 }
133 {
134   "_id" : "قرچک",
135   "value" : 54.74481202576004
136 }
137 {
138   "_id" : "قزوین",
139   "value" : 41.29597490186893
140 }
141 {
142   "_id" : "قم",
143   "value" : 43.728207532391416
144 }

```


▼ (37) { _id : مشهد }	{ 2 fields }	145 {
"_" _id	مشهد	146 "_id" : "مشهد",
1.23 value	49.23186320065565	147 "value" : 49.23186320065565
▼ (38) { _id : ملارد }	{ 2 fields }	148 }
"_" _id	ملارد	149 {
1.23 value	42.97194048938662	150 "_id" : "ملارد",
▼ (39) { _id : نجف آباد }	{ 2 fields }	151 "value" : 42.97194048938662
"_" _id	نجف آباد	152 }
1.23 value	47.142861394210115	153 {
▼ (40) { _id : نیشابور }	{ 2 fields }	154 "_id" : "نجف آباد",
"_" _id	نیشابور	155 "value" : 47.142861394210115
1.23 value	48.82995801792238	156 }
▼ (41) { _id : همدان }	{ 2 fields }	157 {
"_" _id	همدان	158 "_id" : "نیشابور",
1.23 value	43.736670545322376	159 "value" : 48.82995801792238
▼ (42) { _id : ورامین }	{ 2 fields }	160 }
"_" _id	ورامین	161 {
1.23 value	42.81536986098197	162 "_id" : "همدان",
▼ (43) { _id : پاکدشت }	{ 2 fields }	163 "value" : 43.736670545322376
"_" _id	پاکدشت	164 }
1.23 value	51.3831972475595	165 {
▼ (44) { _id : کاشان }	{ 2 fields }	166 "_id" : "ورامین",
"_" _id	کاشان	167 "value" : 42.81536986098197
1.23 value	48.217494459464746	168 }
▼ (45) { _id : کرج }	{ 2 fields }	169 {
"_" _id	کرج	170 "_id" : "پاکدشت",
1.23 value	53.22024724314949	171 "value" : 51.3831972475595
▼ (46) { _id : کرمان }	{ 2 fields }	172 }
"_" _id	کرمان	173 {
1.23 value	46.61853642712546	174 "_id" : "کاشان",
▼ (47) { _id : کرمانشاه }	{ 2 fields }	175 "value" : 48.217494459464746
"_" _id	کرمانشاه	176 }
1.23 value	46.794809852239155	177 {
▼ (48) { _id : گرگان }	{ 2 fields }	178 "_id" : "کرج",
"_" _id	گرگان	179 "value" : 53.22024724314949
1.23 value	46.061132653762606	180 }
▼ (49) { _id : گلستان }	{ 2 fields }	181 {
"_" _id	گلستان	182 "_id" : "کرمان",
1.23 value	44.24247563997863	183 "value" : 46.61853642712546
▼ (50) { _id : یزد }	{ 2 fields }	184 }
"_" _id	یزد	185 {
1.23 value	50.11482544521847	186 "_id" : "کرمانشاه",
		187 "value" : 46.794809852239155
		188 }
		189 {
		190 "_id" : "گرگان",
		191 "value" : 46.061132653762606
		192 }
		193 {
		194 "_id" : "گلستان",
		195 "value" : 44.24247563997863
		196 }
		197 {
		198 "_id" : "یزد",
		199 "value" : 50.11482544521847
		200 }
		201 }

3.278s

زمان اجرا:

بخش دوم: مدل سازی داده ها با گراف / کار با دیتابیس Neo4j

در اینجا ابتدا خلاصه ای از مراحل انجام شده برای پاسخ به سوالات داده شده ارائه می شود:

1. ابتدا نسخه ی 1M دیتاست MovieLens را از آدرس زیر دانلود می کنیم:

<http://files.grouplens.org/datasets/movielens/ml-1m.zip>

این دیتاست شامل 3 فایل `movies.dat`، `users.dat`، `ratings.dat` می باشد:

- فایل `movies.dat` شامل اطلاعات مربوط به 3883 فیلم با فرمت:

MovieID::Title (Year)::Genres

و به صورت زیر می باشد:

```
1::Toy Story (1995)::Animation|Children's|Comedy
2::Jumanji (1995)::Adventure|Children's|Fantasy
3::Grumpier Old Men (1995)::Comedy|Romance
4::Waiting to Exhale (1995)::Comedy|Drama
5::Father of the Bride Part II (1995)::Comedy
6::Heat (1995)::Action|Crime|Thriller
7::Sabrina (1995)::Comedy|Romance
8::Tom and Huck (1995)::Adventure|Children's
9::Sudden Death (1995)::Action
10::GoldenEye (1995)::Action|Adventure|Thriller
11::American President, The (1995)::Comedy|Drama|Romance
```

- فایل `users.dat` شامل اطلاعات مربوط به 6040 کاربر با فرمت:

UserID::Gender::Age::Occupation::Zip-code

و به صورت زیر می باشد:

```
1::F::1::10::48067
2::M::56::16::70072
3::M::25::15::55117
4::M::45::7::02460
5::M::25::20::55455
6::F::50::9::55117
7::M::35::1::06810
8::M::25::12::11413
9::M::25::17::61614
10::F::35::1::95370
11::F::25::1::04093
```

- فایل ratings.dat شامل اطلاعات مربوط به 1000209 مورد Rating با فرمت:

UserID::MovieID::Rating::Timestamp

و به صورت زیر می باشد:

```
1::1193::5::978300760
1::661::3::978302109
1::914::3::978301968
1::3408::4::978300275
1::2355::5::978824291
1::1197::3::978302268
1::1287::5::978302039
1::2804::5::978300719
1::594::4::978302268
1::919::4::978301368
1::595::5::978824268
```

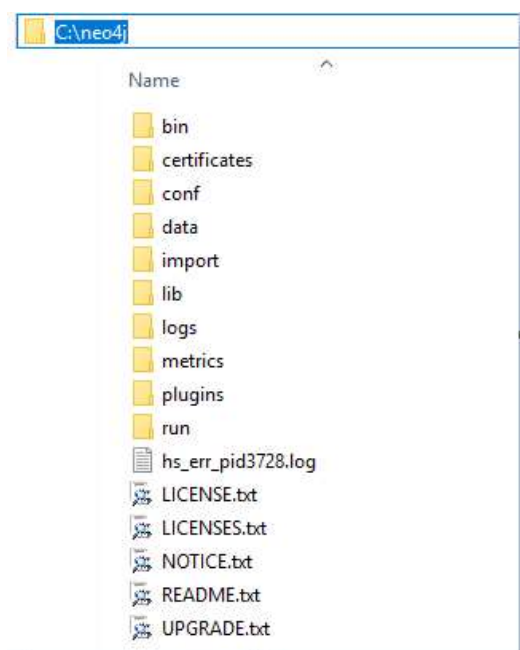
2. سپس نسخه‌ی Enterprise server مربوط به سیستم‌عامل ویندوز دیتاست Neo4j را از آدرس زیر دانلود می‌کنیم:

<https://neo4j.com/download-center/>

Current Releases

Enterprise Server	Community Server	Neo4j Desktop
Neo4j Enterprise Edition 4.0.4 4 May 2020 Release Notes Read More		
OS	Download	
Linux/Mac	Neo4j 4.0.4 (tar) SHA-256	
Windows	Neo4j 4.0.4 (zip) SHA-256	

و سپس تمام موارد داخل این فایل فشرده را در آدرس **C:\neo4j** اکسترکت می‌نماییم. محتوای این آدرس به‌صورت زیر می‌باشد:



سپس آدرس **C:\neo4j\bin** را به Path های سیستم اضافه می کنیم.

و در نهایت، دستور **neo4j install-service** را در cmd اجرا می کنیم تا از این به بعد بتوانیم Neo4j را به عنوان یک windows service داشته باشیم. سپس با دستور **neo4j start**، این سرویس را اجرا می کنیم:

```
C:\Users\Peyman>neo4j start  
Neo4j service started
```


3. سپس دیتاست دانلود شده را از طریق کتابخانه‌ی py2neo نسخه‌ی 4.0.0b1 و به‌صورت زیر در Neo4j درج می‌کنیم:

لازم به ذکر است که با توجه به اینکه دیتاست 1M با دیتاست 100K متفاوت است، به‌همین دلیل مجبور به تغییر بخش‌هایی از کد ارائه شده در آدرس زیر شدیم:

<https://github.com/tkcsdvd/neo4j-movieLens/blob/master/ingestion/ingestion.py>

```
import pandas as pd
from py2neo import Graph, Node

USERNAME = "neo4j"
PASS = "neo4j"

graph = Graph("bolt://localhost:7687", auth = (USERNAME, PASS))
graph.delete_all()

def main():
    print("Create Genre Nodes")
    createGenreNodes()

    print("Step 1 out of 3: loading movie nodes")
    loadMovies()

    print("Step 2 out of 3: loading user nodes")
    loadUsers()

    print("Step 3 out of 3: loading rating relationships")
    loadRatings()

# -----
def createGenreNodes():
    allGenres = ["Action", "Adventure", "Animation", "Children's", "Comedy", "Crime",
                 "Documentary", "Drama", "Fantasy", "Film-Noir", "Horror", "Musical",
                 "Mystery", "Romance", "Sci-Fi", "Thriller", "War", "Western"]

    for genre in allGenres:
        gen = Node("Genre", name=genre)
        graph.create(gen)

# -----
def loadMovies():
    readCSV = pd.read_csv('data/movies.dat', sep=":", engine="python", header=None)
    for row in readCSV.itertuples():
        i = row.Index
        createMovieNodes(row)
        createGenreMovieRelationships(row)
        if i%100==0:
            print(f"{i} Movie nodes created")

# -----
def createMovieNodes(row):
    movieId = row._1
    title = row._2[:-7]
    year = row._2[-5:-1]
    mov = Node("Movie", id=movieId, title=title, year=year)
    graph.create(mov)
```

```

#-----
def createGenreMovieRelationships(row):
    movieId = row._1
    movieGenres = row._3.split("|")

    for movieGenre in movieGenres:
        graph.run('MATCH (g:Genre {name: {genre}}), (m:Movie {id: {movieId}}) CREATE (g)-[:IS_GENRE_OF]->(m)',
            genre=movieGenre, movieId=movieId)

#-----
def loadUsers():
    readCSV = pd.read_csv('data/users.dat', sep=":", engine="python", header=None)
    for row in readCSV.itertuples():
        i = row.Index
        createUserNodes(row)
        if i % 100 == 0:
            print(f"{i} user nodes created")

#-----
def createUserNodes(row):
    userId = "User "+str(row._1)
    gender = row._2
    age = row._3
    occupation = row._4
    zipcode = row._5
    usr = Node("User", id=userId, gender=gender, age=age, occupation=occupation, zipcode=zipcode)
    graph.create(usr)

#-----
def loadRatings():
    readCSV = pd.read_csv('data/ratings.dat', sep=":", engine="python", header=None)
    for row in readCSV.itertuples():
        i = row.Index
        createRatingRelationship(row)
        if i%100==0:
            print(f"{i} Rating relationships created")

#-----
def createRatingRelationship(row):
    ratingData = parseRowRatingRelationships(row)
    graph.run(
        'MATCH (u:User {id: {userId}}), (m:Movie {id: {movieId}}) CREATE (u)-[:RATED { rating: {rating}, timestamp: {timestamp} }]->(m)'
        userId=ratingData[0], movieId=ratingData[1], rating=ratingData[2], timestamp=ratingData[3])

#-----
def parseRowRatingRelationships(row):
    userId = "User "+str(row._1)
    movieId = row._2
    rating = float(row._3)
    timestamp = row._4
    return (userId, movieId, rating, timestamp)

#-----
if __name__ == '__main__':
    main()

```

4. آدرس <http://localhost:7474/browser> را از طریق مرورگر خود باز کرده و کوئری‌های زیر را به‌منظور

آشنایی با زبان Cypher و همچنین بررسی صحت داده‌های درج شده، اجرا می‌کنیم:

- اطلاعات مربوط به یک نود ژانر:

```
$ MATCH (g:Genre) RETURN g LIMIT 1
```

"g"
{"name": "Action"}

- لیست ژانرها و تعداد آنها:

```
$ MATCH (g:Genre) RETURN collect(g.name), count(*)
```

"collect (g.name) "	"count (*) "
["Action", "Adventure", "Animation", "Children's", "Comedy", "Crime", "Docum entary", "Drama", "Fantasy", "Film-Noir", "Horror", "Musical", "Mystery", "Ro mance", "Sci-Fi", "Thriller", "War", "Western"]	18

- اطلاعات مربوط به یک نود فیلم:

```
$ MATCH (m:Movie) RETURN m LIMIT 1
```

"m"
{"title": "Toy Story", "year": "1995", "id": 1}

- تعداد فیلم‌ها:

```
$ MATCH (:Movie) RETURN count(*)
```

"count (*) "
3883

- اطلاعات مربوط به یک نود کاربر:

```
$ MATCH (u:User) RETURN u LIMIT 1
```

Graph	"u"
Table	{"zipcode":"70072","occupation":16,"id":"User 2","gender":"M","age":56}

- تعداد کاربران:

```
$ MATCH (u:User) RETURN count(*)
```

Table	"count (*)"
A	6040

- اطلاعات مربوط به یک یال RATED:

```
$ MATCH (:User)-[r:RATED]-(:Movie) RETURN r LIMIT 1
```

Table	"r"
A	{"rating":3.0,"timestamp":965328072}

- تعداد Rating ها:

```
$ MATCH (:User)-[r:RATED]-(:Movie) RETURN count(*)
```

Table	"count (*)"
A	1000209

سوالات و پرس و جوهای مورد نیاز:

1. چه ژانرهایی در این دیتاست وجود دارد:

Query

```

1 MATCH (g:Genre)
2 WITH collect(g.name) as genres
3 RETURN genres

```

Output

Table

Text

genres

["Action", "Adventure", "Animation", "Children's", "Comedy", "Crime", "Documentary", "Drama", "Fantasy", "Film-Noir", "Horror", "Musical", "Mystery", "Romance", "Sci-Fi", "Thriller", "War", "Western"]

Run time

Started streaming 1 records after 1 ms and completed after 1 ms.

2. تعداد کل فیلم‌ها چقدر است:

Query

```

1 MATCH (:Movie)
2 WITH count(*) as count
3 RETURN count

```

Output

count
3883

Run time

Started streaming 1 records in less than 1 ms and completed in less than 1 ms.

3. تعداد امتیازات داده شده به فیلم Silence of the Lambs در هر رده امتیازی چقدر است:

Query

```
1 MATCH (:User)-[r:RATED]->(m:Movie)
2 WHERE m.title="Silence of the Lambs, The"
3 WITH r.rating as rating, count(*) as count
4 RETURN rating, count
5 ORDER BY rating
```

Output

rating	count
1.0	37
2.0	43
3.0	246
4.0	902
5.0	1350

Run
time

Started streaming 5 records after 20 ms and completed after 20 ms.

4. تعداد امتیازات هر فیلم به صورت نزولی:

Query

```
1 MATCH (:User)-[r:RATED]->(m:Movie)
2 WITH m.id as id, m.title as title, count(*) as count
3 RETURN title, count
4 ORDER BY count DESC
```

Output

title	count
"American Beauty"	3428
"Star Wars: Episode IV - A New Hope"	2991
"Star Wars: Episode V - The Empire Strikes Back"	2990
"Star Wars: Episode VI - Return of the Jedi"	2883
"Jurassic Park"	2672
"Saving Private Ryan"	2653
"Terminator 2: Judgment Day"	2649
"Matrix, The"	2590
"Back to the Future"	2583
"Silence of the Lambs, The"	2578
"Men in Black"	2538

* خروجی کامل این سوال در فولدر Neo4j-Part 2 از فولدر Code قرار گرفته است.

Run
time

Started streaming 3706 records after 3356 ms and completed after 3368 ms, displaying first 1000 rows.

5. کدام ژانر بیشترین فیلم تولید شده را به خود اختصاص داده است:

Query

```
1 MATCH (g:Genre)-[:IS_GENRE_OF]->(m:Movie)
2 WITH g.name as genre, count(*) as count
3 ORDER BY count DESC
4 LIMIT 1
5 RETURN genre
```

Output

genre
"Drama"

Run time

Started streaming 1 records after 11 ms and completed after 12 ms.

کدام ژانر بیشترین امتیاز کاربران را به دست آورده است؟ (بر اساس میانگین امتیاز برای هر ژانر)

Query

```
1 MATCH (g:Genre)-[:IS_GENRE_OF]->(:Movie)<-[:RATED]-(:User)
2 WITH g.name as genre, avg(r.rating) as avgrating
3 ORDER BY avgrating DESC
4 LIMIT 1
5 RETURN genre
```

Output

genre
"Film-Noir"

Run time

Started streaming 1 records after 6161 ms and completed after 6161 ms.

کدام ژانر کمترین امتیاز کاربران را به دست آورده است؟ (بر اساس میانگین امتیاز برای هر ژانر)

Query

```
1 MATCH (g:Genre)-[:IS_GENRE_OF]->(:Movie)<-[:RATED]-(:User)
2 WITH g.name as genre, avg(r.rating) as avgrating
3 ORDER BY avgrating
4 LIMIT 1
5 RETURN genre
```

Output

genre
"Horror"

Run time

Started streaming 1 records after 6084 ms and completed after 6084 ms.

لازم به ذکر است که در قالب یک کوئری نیز می توان به سه سوال ذکر شده پاسخ داد:

Query

```
1 MATCH (g:Genre)-[:IS_GENRE_OF]->(:Movie)
2 WITH g.name as genre, count(*) as count, "max #movies" as QType
3 ORDER BY count DESC
4 LIMIT 1
5 RETURN QType,genre
6 UNION ALL MATCH (g:Genre)-[:IS_GENRE_OF]->(:Movie)<-[r:RATED]-(:User)
7 WITH g.name as genre, avg(r.rating) as avgrating, "max avg rating" as QType
8 ORDER BY avgrating DESC
9 LIMIT 1
10 RETURN QType, genre
11 UNION ALL MATCH (g:Genre)-[:IS_GENRE_OF]->(:Movie)<-[r:RATED]-(:User)
12 WITH g.name as genre, avg(r.rating) as avgrating, "min avg rating" as QType
13 ORDER BY avgrating
14 LIMIT 1
15 RETURN QType, genre
```

Output

	QType	genre
Table	"max #movies"	"Drama"
A	"max avg rating"	"Film-Noir"
Text	"min avg rating"	"Horror"

Run
time

Started streaming 3 records after 18910 ms and completed after 18910 ms.

6. عنوان فیلم‌های تولید شده در سال 2000 را بیابید:

Query

```
1 MATCH (m:Movie)
2 WHERE m.year="2000"
3 WITH m.title as title
4 RETURN title
```

Output

Table	Text	Code	title
			"Supernova"
			"Closer You Get, The"
			"Mission to Mars"
			"Down to You"
			"Isn't She Great?"
			"Scream 3"
			"Gun Shy"
			"Beach, The"
			"Snow Day"
			"Tigger Movie, The"
			"Trois"

* خروجی کامل این سوال در فولدر jNeo4-Part 2 از فولدر Code قرار گرفته است.

Run time Started streaming 156 records in less than 1 ms and completed after 5 ms.

7. با توجه به فیلد occupation در فایل users.dat، برنامه‌نویسان به کدام فیلم بیشترین امتیاز را داده‌اند؟ (کدام فیلم تعداد امتیازی 5 بیشتری دریافت کرده است)

Query

```
1 MATCH (u:User)-[r:RATED]->(m:Movie)
2 WHERE u.occupation=12 and r.rating=5
3 WITH m.id as id, m.title as title, count(*) as count
4 ORDER BY count DESC
5 LIMIT 1
6 RETURN title
```

Output

Table	title
	"Star Wars: Episode IV - A New Hope"

Run time Started streaming 1 records after 348 ms and completed after 348 ms.

8. پنج فیلم محبوب رده‌ی سنی 18 تا 34 سال را به‌دست آورید. (بر اساس میانگین امتیاز)

Query

```
1 MATCH (u:User)-[r:RATED]->(m:Movie)
2 WHERE u.age=18 or u.age=25
3 WITH m.id as id, m.title as title, avg(r.rating) as avgrating
4 ORDER BY avgrating DESC
5 LIMIT 5
6 RETURN title
```

Output

	title
Table	"Callejón de los milagros, El"
Text	"Gate of Heavenly Peace, The"
	"Baby, The"
Code	"Dingo"
	"Black Sunday (La Maschera Del Demonio)"

* با توجه به اینکه در اینجا تعداد زیادی فیلم، میانگین امتیاز برابر 5 دارند، به‌همین دلیل اگر بخواهیم فقط 5 مورد از آنهایی که بیشترین میانگین امتیاز (در اینجا، 5) داشته‌اند را انتخاب کنیم، آنگاه این 5 مورد می‌تواند متفاوت باشد.

Run
time

Started streaming 5 records after 3125 ms and completed after 3125 ms.

9. بیست فیلم محبوب که حداقل 100 کاربر به آنها رای داده باشند را بدست آورید. (بر اساس

میانگین امتیاز)

Query

```
1 MATCH (u:User)-[r:RATED]->(m:Movie)
2 WITH m.id as id,m.title as title, count(*) as count, avg(r.rating) as avgrati
3 WHERE count>=100
4 RETURN title
5 ORDER BY avgrating DESC
6 LIMIT 20
```

Output

Table	title
	"Seven Samurai (The Magnificent Seven) (Shichinin no samurai)"
	"Shawshank Redemption, The"
	"Godfather, The"
	"Close Shave, A"
	"Usual Suspects, The"
	"Schindler's List"
	"Wrong Trousers, The"
	"Sunset Blvd. (a.k.a. Sunset Boulevard)"
	"Raiders of the Lost Ark"
	"Rear Window"
	"Paths of Glory"
	"Star Wars: Episode IV - A New Hope"
	"Third Man, The"
	"Dr. Strangelove or: How I Learned to Stop Worrying and Love the Bomb"
	"Wallace & Gromit: The Best of Aardman Animation"
	"To Kill a Mockingbird"
	"Double Indemnity"
	"Casablanca"
	"Sixth Sense, The"
	"Yojimbo"

Run
time

Started streaming 20 records after 3780 ms and completed after 3780 ms.

بخش سوم: دیتابیس‌های سطرگسترده / کار با HBase

گام اول:

به‌منظور پاسخ به کوئری‌های داده شده، تنها از دو جدول به‌صورت زیر استفاده خواهیم کرد:

1. جدول ژانر:

این جدول، شامل اطلاعات فیلم‌های مربوط به هر ژانر است.

در اینجا کلید سطر، نام ژانر است و یک ColumnFamily به نام isgenreof ایجاد شده است که هر ستون آن، اطلاعات یک فیلم مربوط به آن ژانر را نشان می‌دهد. این جدول به‌صورت زیر می‌باشد:

Row id	Column family		
genre name	isgenreof		
	isgenreof:movieId1	isgenreof:movieId2	isgenreof:movieId3
Action	movie title 1	movie title 2	movie title 3
Adventure	movie title 4	movie title 5	-

همچنین در شکل زیر که سه سطر نمونه از این جدول را نشان می‌دهد، توضیحات تکمیلی ارائه شده است:

ColumnFamily isgenreof		
ژانر		
Action		
isgenreof: 1129	isgenreof: 6	isgenreof: 748
Escape from New York	Heat	Arrival, The
Adventure		
isgenreof: 2	isgenreof: 1129	
Jumanji	Escape from New York	
Animation		
isgenreof: 1		
Toy Story		

ما از این جدول برای پاسخ به سوال 1 استفاده می‌نماییم.

2- جدول فیلم:

این جدول، شامل اطلاعات مربوط به هر کدام از فیلم‌ها و rating‌های داده شده به آنها می‌باشد.

در اینجا، کلید سطر، شناسه‌ی فیلم است و دو ColumnFamily به شرح زیر ایجاد شده است:

1. data که اطلاعات اصلی یک فیلم را ذخیره می‌کند، یعنی شامل دو ستون data:title و data:year می‌باشد.

2. rating که اطلاعات مربوط به rating‌های داده شده به فیلم مربوطه را ذخیره می‌کند، یعنی شامل ستون‌هایی به صورت rating:userId می‌باشد. هر ستون rating:userId شامل مقادیر gender::age::rating_score از آن کاربر خواهد بود.

این جدول به صورت زیر می‌باشد:

Row_id	Column family		Column family		
movieId	data		rating		
	data::title	data:year	rating:userId	rating:userId2	...
movieId1	title name	year value	gender::age::rating_value

همچنین در شکل زیر سه سطر نمونه از این جدول را نشان می‌دهد:

1	rating: 6016	rating: 6013	data: year	rating: 6035	data: title	genre: genre2	genre: genre3
	F::25::4	M::18::5	1995	M::58::4	Toy Story	Children's	Comedy
1129	data: title	rating: 6003	rating: 6007	rating: 5947	rating: 6030	rating: 6036	data: year
	Escape from New York	M::45::14	F::18::3	F::25::2	M::58::4	M::35::3	1981
1316	data: year	data: title	rating: 3952	genre: genre1			
	1996	Anna	F::45::4	Drama			

ما از این جدول برای پاسخ به سوالات 2 تا 5 استفاده می‌نماییم.

گام دوم:

نمونه‌هایی از داده‌های دستی وارد شده به دو جدول ایجاد شده، در صفحات قبل ارائه شد. در اینجا بر اساس داده‌های وارد شده، به سوالات زیر پاسخ می‌دهیم:

1. فهرست فیلم‌های یک ژانر:

با توجه به اینکه گفته شده فهرست فیلم‌های یک ژانر تولید شود، به همین دلیل ما در اینجا به عنوان نمونه، فیلم‌های ژانر Thriller را تولید خواهیم کرد.

دستورات نوشته شده:

```
import happybase

connection = happybase.Connection(host='quickstart.cloudera',port=9090)
connection.open()

table = connection.table('Genre')
for key,value in table.scan(row_start='Thriller',row_stop='Thriller',columns=['isgenreof']):
    print(value.values())
```

خروجی:

```
sh-4.1# python3 code/hbase.py
dict_values([b'Heat', b'Arrival, The', b'Escape from New York'])
```

2. تعداد ژانرهای هر فیلم:

دستورات نوشته شده:

```
import happybase

connection = happybase.Connection(host='quickstart.cloudera',port=9090)
connection.open()

table = connection.table('Movie')
for key,value in table.scan(columns=['data:title','genre']):
    print(value[b'data:title'],len(value)-1)
```

خروجی:

```
sh-4.1# python3 code/hbase.py
b'Toy Story' 3
b'Escape from New York' 4
b'Anna' 1
b'Jumanji' 3
b'Heat' 3
b'Arrival, The' 3
```

3. کدام فیلم بیشترین تعداد امتیازدهنده را داشته است:

دستورات نوشته شده:

```
import happybase
import numpy as np

connection = happybase.Connection(host='quickstart.cloudera',port=9090)
connection.open()

table = connection.table('Movie')

title_list, ratingcount_list = list(), list()
for key,value in table.scan(columns=['data:title','rating']):
    title_list.append(value[b'data:title'])
    ratingcount_list.append(len(value)-1)

title = title_list[np.argmax(ratingcount_list)]
print(title)
```

خروجی:

```
sh-4.1# python3 code/hbase.py
b'Arrival, The'
```

4. سه فیلم محبوب تر بازه سنی 45 سال به بالا:

دستورات نوشته شده:

```
import happybase
import numpy as np

connection = happybase.Connection(host='quickstart.cloudera',port=9090)
connection.open()

table = connection.table('Movie')

def avg_per_movie(value):
    avg = np.mean([int(user_rating[7:]) for user_rating in value.values()])
    return avg

dic = dict()
for key,value in table.scan(filter=" ValueFilter (='regexstring:^[MF]::45::[1-5]$') OR"
                             " ValueFilter (='regexstring:^[MF]::50::[1-5]$') OR"
                             " ValueFilter (='regexstring:^[MF]::56::[1-5]$')"):
    dic[key]=avg_per_movie(value)

best_movies = [int(k) for k,v in sorted(dic.items(), key=lambda item: -item[1])]
print(best_movies[:3])
```

خروجی:

```
sh-4.1# python3 code/hbase.py
[1316, 1, 1129]
```

یعنی فیلم‌هایی با سه شناسه فیلم 1316، 1 و 1129، به ترتیب سه فیلم محبوب بازه سنی 45 سال به بالا هستند.

5. میانگین امتیاز خانم‌ها و آقایان در فیلم‌های سال 1995:

دستورات نوشته شده برای یافتن میانگین امتیاز آقایان در فیلم‌های سال 1995:

```
import happybase
import numpy as np

connection = happybase.Connection(host='quickstart.cloudera',port=9090)
connection.open()

table = connection.table('Movie')

avglist = list()
for key,value in table.scan(filter="SingleColumnValueFilter ('data','year',, 'binary:1995') AND "
                             "ValueFilter (='regexstring: ^M::[0-9]{2}::[1-5]$')"):
    avglist += list(value.values())

avg = np.mean([int(value[7:]) for value in avglist])
print(avg)
```

خروجی:

```
sh-4.1# python3 code/hbase.py
3.375
```

دستورات نوشته شده برای یافتن میانگین امتیاز خانم‌ها در فیلم‌های سال 1995:

```
import happybase
import numpy as np

connection = happybase.Connection(host='quickstart.cloudera',port=9090)
connection.open()

table = connection.table('Movie')

avglist = list()
for key,value in table.scan(filter="SingleColumnValueFilter ('data','year',, 'binary:1995') AND "
                             "ValueFilter (='regexstring: ^F::[0-9]{2}::[1-5]$')"):
    avglist += list(value.values())

avg = np.mean([int(value[7:]) for value in avglist])
print(avg)
```

خروجی:

```
sh-4.1# python3 code/hbase.py
4.0
```