https://learning.oreilly.com/library/view/mongodb-in-action/9781617291609/?ar (Links to an external site.)

Chapter 4. Document-oriented data

Chapter 5. Constructing queries

Chapter 6. Aggregation

Note: If you cannot access the chapters, enter your neu email as @northeastern.edu instead of @husky.neu.edu
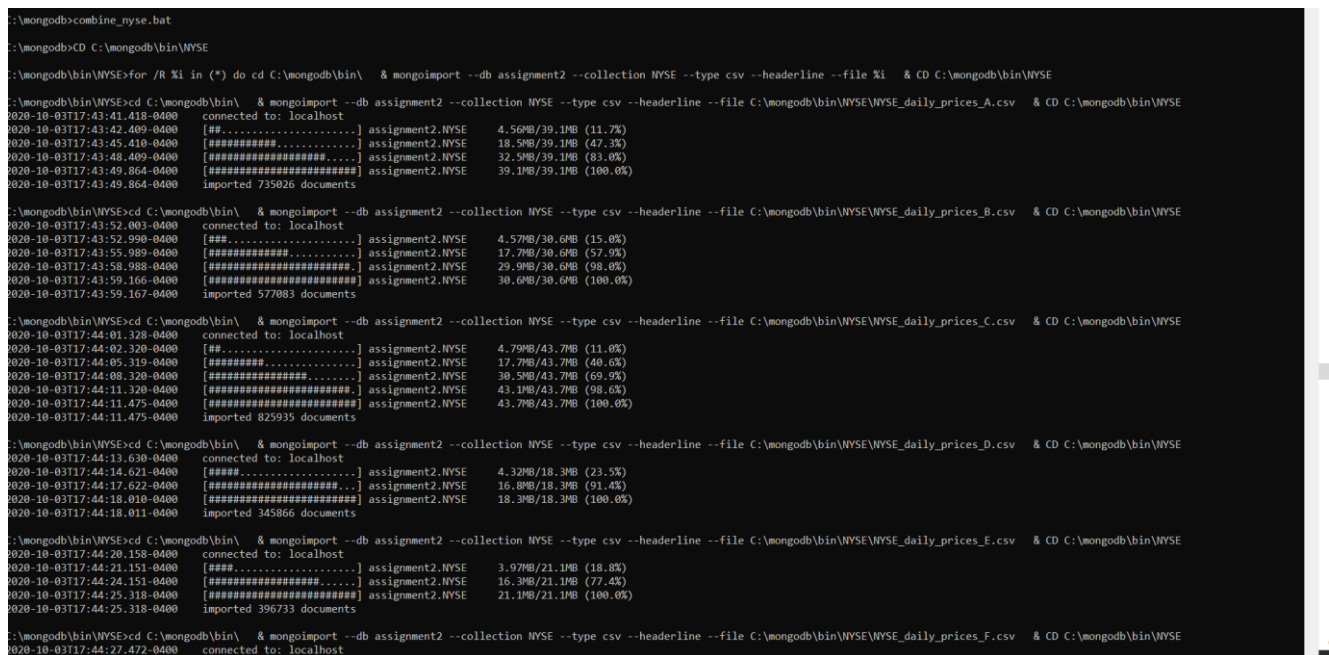
**PART 2 - PROGRAMMING ASSIGNMENT**

Write a .bat/.sh to import the entire NYSE dataset (stocks A to Z) into MongoDB.

NYSE Dataset Link: http://msis.neu.edu/nyse/



```
C:\mongodb\combine_nyse.bat - Notepad++
File Edit Search View Encoding Language Settings Tools Macro Run Plugins Window ?

1  CD C:\mongodb\bin\NYSE
2  for /r %%i in (*) do cd C:\mongodb\bin\ & mongoimport --db assignment2 --collection NYSE --type csv --headerline --file %%i & CD C:\mongodb\bin\NYSE
```

PART 3.1. Use the NYSE database to find the average price of stock_price_high values for each stock using MapReduce.

```
var map = function() {emit(this.stock_symbol, this.stock_price_high);};
var reduce = function(sym, high){return Array.avg(high);  };

db.NYSE.mapReduce(map, reduce, {out : "avg_high_price"}).find();
```

avg_high_price    🕐 146 sec.

| Key | Value | Type |
|---|---|---|
| ∨ 💿 (1) NaN | { 2 fields } | Object |
|   📄 _id | NaN | Double |
|   📄 value | 14.3482400216403 | Double |
| ∨ 💿 (2) AA | { 2 fields } | Object |
|   📄 _id | AA | String |
|   📄 value | 64.2280328504867 | Double |
| ∨ 💿 (3) AAI | { 2 fields } | Object |
|   📄 _id | AAI | String |
|   📄 value | 21.8955523294756 | Double |
| ∨ 💿 (4) AAN | { 2 fields } | Object |
|   📄 _id | AAN | String |
|   📄 value | 7.84126067143493 | Double |
| > 💿 (5) AAP | { 2 fields } | Object |
| > 💿 (6) AAR | { 2 fields } | Object |
| > 💿 (7) AAV | { 2 fields } | Object |
| > 💿 (8) AB | { 2 fields } | Object |
| > 💿 (9) ABA | { 2 fields } | Object |
| > 💿 (10) ABB | { 2 fields } | Object |
| > 💿 (11) ABC | { 2 fields } | Object |
| > 💿 (12) ABD | { 2 fields } | Object |
| > 💿 (13) ABG | { 2 fields } | Object |
| > 💿 (14) ABK | { 2 fields } | Object |
| > 💿 (15) ABM | { 2 fields } | Object |
| > 💿 (16) ABR | { 2 fields } | Object |
| > 💿 (17) ABT | { 2 fields } | Object |
| > 💿 (18) ABV | { 2 fields } | Object |
| > 💿 (19) ABVT | { 2 fields } | Object |
| > 💿 (20) ABX | { 2 fields } | Object |
| > 💿 (21) ACC | { 2 fields } | Object |
| > 💿 (22) ACE | { 2 fields } | Object |
| > 💿 (23) ACF | { 2 fields } | Object |
| > 💿 (24) ACG | { 2 fields } | Object |

PART 3.2. Part 3.1 result will not be correct as AVERAGE is a commutative operation but nor associative. Use a FINALIZER to find the correct average. (Hint: pass sum and count from the reducer) (https://docs.mongodb.com/manual/reference/method/db.collection.mapReduce/index.html)

New Connection    localhost:27017    assignment2

```
var map = function(){
    emit (this.stock_symbol , {high_price : this.stock_price_high, count: 1});};

var reduce = function(sym, values){
    reduce_output = {high_price:0, count:0};
    values.forEach(function(value){
        reduce_output.high_price += value.high_price;
        reduce_output.count += value.count;
    });
    return reduce_output;
}

var finalizefn = function(sym, reduceVal){
    reduceVal.avg = reduceVal.high_price / reduceVal.count;
        return reduceVal;
    }

db.NYSE.mapReduce(map, reduce, {out : "count_total", finalize:finalizefn }).find();
```

count_total    🕐 98.2 sec.                                         ◄  0    50  ►

| Key | Value | Type |
|---|---|---|
| > 💿 (1) NaN | { 2 fields } | Object |
| ∨ 💿 (2) AA | { 2 fields } | Object |
|   📄 _id | AA | String |
|   ∨ 💿 value | { 3 fields } | Object |
|     📄 high_price | 635234.290000002 | Double |
|     📄 count | 12109.0 | Double |
|     📄 avg | 52.4596820546702 | Double |
| ∨ 💿 (3) AAI | { 2 fields } | Object |
|   📄 _id | AAI | String |
|   ∨ 💿 value | { 3 fields } | Object |
|     📄 high_price | 41369.0500000004 | Double |
|     📄 count | 3933.0 | Double |
|     📄 avg | 10.5184464785152 | Double |
| ∨ 💿 (4) AAN | { 2 fields } | Object |
|   📄 _id | AAN | String |
|   ∨ 💿 value | { 3 fields } | Object |
|     📄 high_price | 83717.149999999 | Double |
|     📄 count | 4218.0 | Double |
|     📄 avg | 19.8475936462776 | Double |
| > 💿 (5) AAP | { 2 fields } | Object |
| > 💿 (6) AAP | { 2 fields } | Object |

## PART 4 - PROGRAMMING ASSIGNMENT

PART 4. Calculate the average stock price of each price of all stocks using $avg aggregation.

https://docs.mongodb.com/manual/reference/operator/aggregation/avg/ (Links to an external site.)

New Connection    localhost:27017    assignment2

```
db.NYSE.aggregate([
    { "$group": {
        "_id": "$stock_symbol",
        "avgQuantity": { $avg: "$stock_price_high" }
    }},
    { "$out": "avg" }
])
db.getCollection('avg').find({})
```

NYSE    21 sec.

Fetched 0 record(s) in 0ms

avg    0.003 sec.

| Key | Value | Type |
| --- | --- | --- |
| ∨ (1) ZMH | { 2 fields } | Object |
| _id | ZMH | String |
| avgQuantity | 61.7007966260544 | Double |
| ∨ (2) ZQK | { 2 fields } | Object |
| _id | ZQK | String |
| avgQuantity | 17.905243902439 | Double |
| ∨ (3) ZAP | { 2 fields } | Object |
| _id | ZAP | String |
| avgQuantity | 11.207132887899 | Double |
| > (4) ZF | { 2 fields } | Object |
| > (5) ZNT | { 2 fields } | Object |
| > (6) ZLC | { 2 fields } | Object |
| > (7) ZEP | { 2 fields } | Object |
| > (8) ZZ | { 2 fields } | Object |

Import the Movielens dataset into MongoDB. Refer to README about file contents and headings.

https://grouplens.org/datasets/movielens/1m/ (Links to an external site.)   [you may replace :: in the dateset with comma or tab to import]

- Find the number Females and Males from the users collection using MapReduce. Do the same thing using count() to compare the results.

```
var map = function(){emit(this.Gender,1);};
var reduce = function(gender,count){
    var sum = Array.sum(count);
    return sum
};
db.users.mapReduce(map, reduce, {out:"gender_count"}) find();

|
```

gender_count    0.284 sec.

| Key | Value | Type |
|---|---|---|
| ✓ (1) F | { 2 fields } | Object |
| _id | F | String |
| value | 1709.0 | Double |
| ✓ (2) M | { 2 fields } | Object |
| _id | M | String |
| value | 4331.0 | Double |

```
db.getCollection('users').find({Gender:"M"}).count()
db.getCollection('users').find({Gender:"F"}).count()
```

0.014 sec.

4331

< 

0.013 sec.

1709

- Find the number of Movies per year using MapReduce

```
var map = function(){
    this.year = this.Title.slice(-5, -1);
    if (isNaN(this.year))
    {
        emit("N / A", 1)
    }
    else
    {
        emit(this.year, 1)
    }
}
var reduce = function(year, count){
    var sum = Array.sum(count);
    return sum;
}
db.movies.mapReduce(map, reduce, {out:"every_movie_year"}).find();
```

every_movie_year   🕐 0.185 sec.

| Key | Value | Type |
|---|---|---|
| ∨ ▣ (1) 1919 | { 2 fields } | Object |
|    _id | 1919 | String |
|    value | 3.0 | Double |
| ∨ ▣ (2) 1920 | { 2 fields } | Object |
|    _id | 1920 | String |
|    value | 2.0 | Double |
| ∨ ▣ (3) 1921 | { 2 fields } | Object |
|    _id | 1921 | String |
|    value | 1.0 | Double |
| > ▣ (4) 1922 | { 2 fields } | Object |
| > ▣ (5) 1923 | { 2 fields } | Object |
| > ▣ (6) 1925 | { 2 fields } | Object |
| > ▣ (7) 1926 | { 2 fields } | Object |
| > ▣ (8) 1927 | { 2 fields } | Object |
| > ▣ (9) 1928 | { 2 fields } | Object |

- Find the number of Movies per rating using MapReduce

```
var map = function(){emit(this.Rating,1);};
var reduce = function(rating,count){
    var sum = Array.sum(count);
    return sum
};
db.ratings.mapReduce(map, reduce, {out:"rating"}).find();
```

rating   🕐 15.3 sec.

| Key | Value | Type |
|---|---|---|
| ∨ ▣ (1) 1.0 | { 2 fields } | Object |
|    _id | 1.0 | Double |
|    value | 56174.0 | Double |
| ∨ ▣ (2) 2.0 | { 2 fields } | Object |
|    _id | 2.0 | Double |
|    value | 107557.0 | Double |
| ∨ ▣ (3) 3.0 | { 2 fields } | Object |
|    _id | 3.0 | Double |
|    value | 261197.0 | Double |
| ∨ ▣ (4) 4.0 | { 2 fields } | Object |
|    _id | 4.0 | Double |
|    value | 348971.0 | Double |
| ∨ ▣ (5) 5.0 | { 2 fields } | Object |
|    _id | 5.0 | Double |
|    value | 226310.0 | Double |