

Aggregation Pipeline

Aldona Biewska

10 maja 2017

- ▶ *Zgony w Stanach Zjednoczonych w 2014 roku.*
- ▶ Plik DeathRecords.csv obrobiony przy pomocy skryptu w R *prepare.R*.
- ▶ Wybrane pliki:

Plik	Liczba dokumentów	Czas importu
DeathRecords__prepare.csv	2631171	4 min 40 s
EntityAxisCondition.csv	8052877	7min 41s
Icd10Code.csv	12131	3s

Przykładowy dokument z kolekcji *deaths*

```
1  {
2    "_id" : ObjectId("58fcdd07bd26e96c4d2d9129"),
3    "Id" : 1,
4    "Sex" : "M",
5    "Age" : 87,
6    "MonthOfDeath" : 1,
7    "Icd10Code" : "I64",
8    "AgeType" : "Years",
9    "Education" : "9 - 12th grade, no diploma",
10   "MaritalStatus" : "Married",
11   "DayOfWeekOfDeath" : "Wednesday",
12   "Race" : "White",
13   "MannerOfDeath" : "Natural",
14   "ActivityCode" : "Not applicable"
15 }
```

Przykładowy dokument z kolekcji *conditions* i *icd10code*

conditions - odzwierciedlenie struktury aktu zgonu:

```
1  {
2    "_id" : ObjectId("58fcdebfbfd26e96c4d55c198"),
3    "Id" : 1,
4    "DeathRecordId" : 1,
5    "Part" : 1,
6    "Line" : 1,
7    "Sequence" : 1,
8    "Icd10Code" : "I64"
9  }
```

icd10code:

```
1  {
2    "_id" : ObjectId("58fce11abd26e96c4dd0c1d1"),
3    "Code" : "A00",
4    "Description" : "Cholera"
5  }
```

Agregacje w JS dostępne są *tutaj*.

Dodatkowo wszystkie agregacje zostały napisane i uruchomione w R, dzięki czemu łatwiej było wykonać wszystkie wykresy oraz tabele *mongo.R*

Agregacja 1: Porównanie samobójstw w danym przedziale wiekowym

Czas wykonania: 4s 370ms

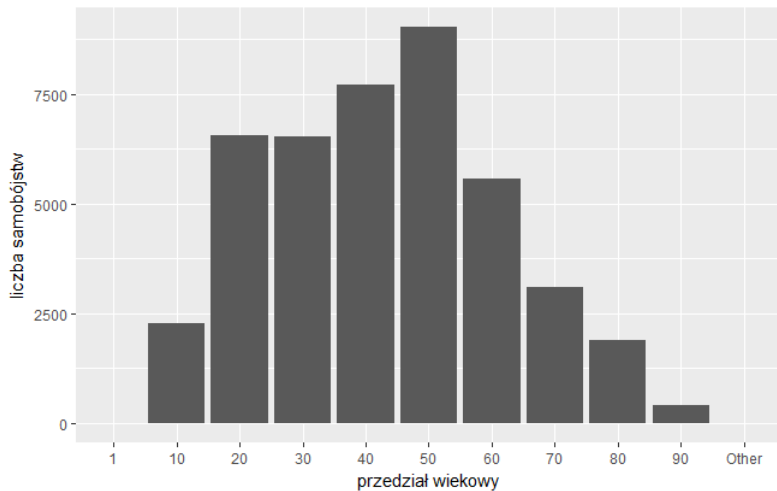
```
1 db.deaths.aggregate([
2   {
3     $match:{
4       "MannerOfDeath":"Suicide",
5       "AgeType":"Years"
6     }
7   },
8   {
9     $bucket:{
10      groupBy:"$Age",
11      boundaries:[1,10,20,30,
12        40,50,60,70,80,90,100],
13      default:"Other"
14    }
15  }
16 ])
```

Wynik:

```
1  [
2    { "_id":1,"count":3 },
3    { "_id":10,"count":2270},
4    { "_id":20,"count":6578},
5    { "_id":30,"count":6526},
6    { "_id":40,"count":7712},
7    { "_id":50,"count":9032},
8    { "_id":60,"count":5588},
9    { "_id":70,"count":3115},
10   { "_id":80,"count":1898},
11   { "_id":90,"count":403},
12   { "_id":"Other","count":7}
13 ]
```

Przy pomocy *\$bucket* i *\$boundaries* określane są przedziały, np.: dla 1 przedziałem jest [1,9).

Agregacja 1: Porównanie samobójstw w danym przedziale wiekowym - Wykres



Agregacja 2: Wpływ edukacji i stanu cywilnego na żywotność kobiet i mężczyzn - krok 1

Wybranie typu wieku (w tym przypadku lata) oraz przedziału od 35 do 100.

```
1  var stage1 = {  
2    $match: {  
3      "AgeType": "Years",  
4      "Education":{$ne: "NA"},  
5      $and: [{ "Age" : { $gt : 35 }}, { "Age" : { $lt : 100 }}]  
6    }  
7  };
```


Agregacja 2: Wpływ edukacji i stanu cywilnego na żywotność kobiet i mężczyzn - krok 2

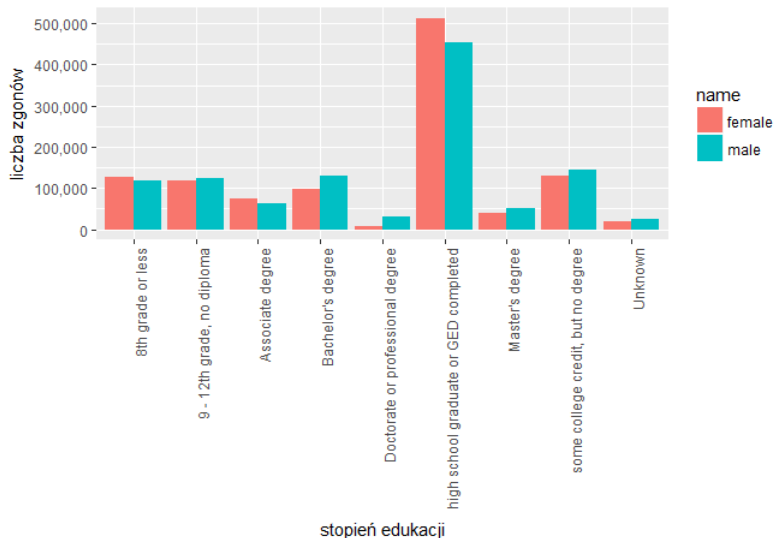
```
1  var stage2 = {
2    $facet: {
3      Education: [
4        { $group: {
5          _id:{$sex: "$Sex",edu: "$Education"},
6          count: {$sum: 1}}},
7        { $sort: {count: -1} },
8        { $group: {
9          _id: "$_id.sex",
10         education: {
11           $push:{range: "$_id.edu" , total:"$count"}}}}}
12      ],
13      Marriage: [
14        { $group: {
15          _id:{$sex: "$Sex", status: "$MaritalStatus" },
16          count: {$sum: 1}} },
17        { $sort: {count: -1} },
18        { $group: {
19          _id: "$_id.sex",
20          marriage: {
21            $push: {status: "$_id.status" , total:"$count"}}}}}
22      ]
23    }
24  };
```

Agregacja 2: Wpływ edukacji i stanu cywilnego na żywotność kobiet i mężczyzn

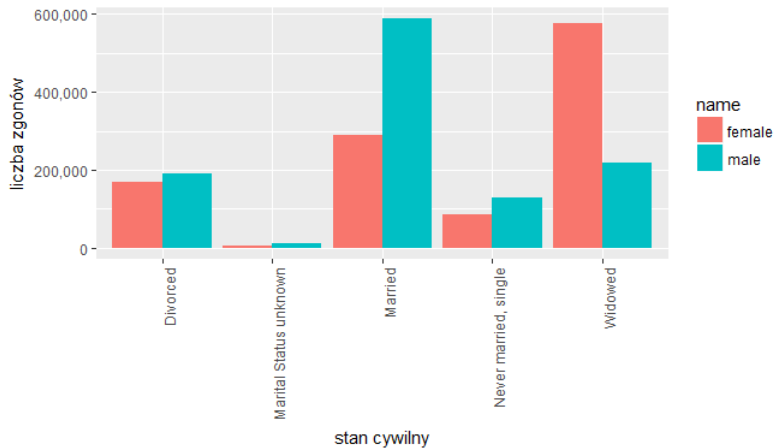
Czas wykonania: 11s 389ms Całość:

```
1 db.deaths.aggregate( [  
2   stage1,  
3   stage2  
4 ] );
```

Agregacja 2: Wpływ edukacji i stanu cywilnego na żywotność kobiet i mężczyzn - Wykres



Agregacja 2: Wpływ edukacji i stanu cywilnego na żywotność kobiet i mężczyzn - Wykres



Agregacja 3: Porównanie ilościowe i procentowe zgonów w wyniku zabójstwa w wybranych rasach - agregacja pomocnicza krok 1

Zliczono liczbę wszystkich zgonów według wybranych ras, przy czym określone rasy szczegółowe przypisano do rasy *Yellow*. Rasa *Black* i *White* była już określona w kolekcji. Rasy dodawane są do nowego pola *Race_new*.

```
1  var stage1 = {
2    $project: {
3      "Id": 1,
4      "MannerOfDeath": 1,
5      "Race_new": {
6        $switch: {
7          branches: [
8            { "case": { $eq: ["$Race", "White"] }, then: "White" },
9            { "case": { $eq: ["$Race", "Black"] }, then: "Black" },
10           { "case": { $or: [
11             { $eq: ["$Race", "Chinese"] },
12             { $eq: ["$Race", "Japanese"] },
13             { $eq: ["$Race", "Asian Indian"] },
14             { $eq: ["$Race", "Korean"] },
15             { $eq: ["$Race", "Other Asian or Pacific Islander"] },
16             { $eq: ["$Race", "American Indian (...)] },
17             { $eq: ["$Race", "Vietnamese"] },
18             { $eq: ["$Race", "Guamanian"] } ] } },
19           then: "Yellow" } ],
20    default: "Did not match" } } };
```

Agregacja 3: Porównanie ilościowe i procentowe zgonów w wyniku zabójstwa w wybranych rasach - agregacja pomocnicza krok 2,3,4

Grupowanie zgonów według ras, odpowiednie wyświetlenie wyników:

```
1  var stage2={
2    $match: { "Race_new": {$ne: "Did not match"}}
3  };
4  var stage3={
5    $group: { "_id": { "race": "$Race_new"}, "count": {$sum: 1}}
6  };
7  var stage4={
8    $project: { "_id": 0, "race": "$_id.race", "count": $count}
9  };
```

Agregacja 3: Porównanie ilościowe i procentowe zgonów w wyniku zabójstwa w wybranych rasach - agregacja pomocnicza

Całość:

```
1  var count =  
2    db.deaths.aggregate([  
3      stage1,  
4      stage2,  
5      stage3,  
6      stage4  
7    ]);
```

Otrzymany wynik:

```
1  { "race" : "Yellow",  
2    "count" : 56205 }  
3  { "race" : "Black",  
4    "count" : 309504 }  
5  { "race" : "White",  
6    "count" : 2241510 }
```

Przypisanie wyników zmiennym:

```
1  var white,yellow,black;  
2  count.forEach(function(record) {  
3    if (record.race=="White") white=record.count;  
4    if (record.race=="Black") black=record.count;  
5    if (record.race=="Yellow") yellow=record.count;  
6  });
```

Agregacja 3: Porównanie ilościowe i procentowe zgonów w wyniku zabójstwa w wybranych rasach - agregacja właściwa krok 1,2,3

Krok 1 taki jak w agregacji pomocniczej.

Krok 2,3 wybieranie rodzaju śmierci jako zabójstwo i grupowanie według ras.

```
1  var stage2_2 = {  
2    $match: {  
3      "MannerOfDeath": "Homicide",  
4      "Race_new": { $ne: "Did not match" } }  
5  };
```

```
1  var stage3_1 = {  
2    $group: { _id: { race: "$Race_new" }, count: { $sum: 1 } }  
3  };
```


Agregacja 3: Porównanie ilościowe i procentowe zgonów w wyniku zabójstwa w wybranych rasach - agregacja właściwa krok 4

Obliczenie jaki procent stanowią zgony w wyniku zabójstwa wobec wszystkich zgonów:

```
1  var stage4_2 = {
2    $project: {
3      "percentage of all deaths":
4        {$let: {
5          vars: { "white": white, "black":black,"yellow":yellow },
6          in: { $switch: {
7            branches: [
8              {case: { $eq: [ "$_id.race","White" ] },
9                then: {
10                  "$multiply":[{"$divide":[100,"$$white"]},"$count"]}},
11              {case: { $eq: [ "$_id.race","Black" ] },
12                then: {
13                  "$multiply":[{"$divide":[100,"$$black"]},"$count"]}},
14              {case: { $eq: [ "$_id.race","Yellow" ] },
15                then: {
16                  "$multiply":[{"$divide":[100,"$$yellow"]},"$count"]}},
17            ],
18            default: "0"}}}},
19    "_id":0,"count":"$count","race":"$_id.race"};
```

Agregacja 3: Porównanie ilościowe i procentowe zgonów w wyniku zabójstwa w wybranych rasach - Wyniki

Czas wykonania: 3s 602ms Całość:

```
1 db.deaths.aggregate([
2   stage1,
3   stage2_2,
4   stage3_2,
5   stage4_2
6 ]);
```

percentage homicide - all deaths	count	race
0.800640512409928	450	Yellow
2.6277527915632755	8133	Black
0.3606051277933179	8083	White

Agregacja 4: Najczęstsze czynniki pośrednie przyczyniające się do głównej przyczyny śmierci

Główna przyczyna: I469 - Cardiac arrest, cause unspecified (Zatrzymanie krążenia)
Dla każdego zgonu tworzę listę przyczyn (2-3); następnie wyszukuję rekordy, w których na liście jako główna przyczyna znajduje się kod *I469* (5-6); wybieram pośrednie przyczyny (8-10) i grupuję (11); łączę z tabelą *icd10*, aby wyodrębnić nazwy chorób (14-18); odpowiednio wyświetlam (19-25).

```
1 db.conditions.aggregate([
2   { $group: {
3     _id : "$DeathRecordId",
4     Other_conditions: { $push: "$$ROOT" } }},
5   { $match: {
6     "Other_conditions":{"$elemMatch":{"Part":1,"Line":1,"Icd10Code":"I469"}}}},
7   { $unwind : "$Other_conditions" },
8   { $match: {
9     "Other_conditions.Icd10Code":{"$ne:"I469"},
10    "Other_conditions.Part":2}},
11   { $group: { _id: "$Other_conditions.Icd10Code", count: { $sum: 1 } }},
12   { $sort: { count: -1 } },
13   { $limit : 5 },
14   { $lookup: {
15     from:"icd10",
16     localField:"_id",
17     foreignField:"Code",
18     as:"Icd10Description"}},
19   { $unwind : "$Icd10Description" },
20   { $project : {
21     "_id":0,
22     "code":"$_id",
23     "description":"$Icd10Description.Description",
24     "count_of_cases":"$count"}}
25 ],{allowDiskUse: true});
```

Agregacja 4: Najczęstsze czynniki pośrednie przyczyniające się do głównej przyczyny śmierci - Wyniki po angielsku

code	description	count
I10	Essential (primary) hypertension	28771
F179	Mental and behavioural disorders due to use of tobacco: Unspecified mental and behavioural disorder	25154
E149	Unspecified diabetes mellitus: Without complications	14413
J449	Chronic obstructive pulmonary disease, unspecified	11588
F03	Unspecified dementia	11165

Agregacja 4: Najczęstsze czynniki pośrednie przyczyniające się do głównej przyczyny śmierci - Wyniki po polsku

code	description	count
I10	Nadciśnienie samoistne (pierwotne)	28771
F179	Zaburzenia psychiczne i zaburzenia zachowania spowodowane paleniem tytoniu (zaburzenia psychiczne i zaburzenia zachowania, nieokreślone)	25154
E149	Cukrzyca nieokreślona (bez powikłań)	14413
J449	Nieokreślona przewlekła zaporowa choroba płuc	11588
F03	Otępienie bliżej nieokreślone	11165