



Politechnika Krakowska
Wydział Informatyki i Telekomunikacji

Projekt z przedmiotu
Metody i narzędzia analizy
dużych zbiorów danych

Projekt nr 1

Kierunek: Informatyka

Stopień studiów: II stopnia

Specjalizacja: Data Science

Wykonali:

Marek Dalida

Maciej Gicala

Rafał Gęgotek

1) Cel projektu	3
2) Zbiór danych	3
3) Implementacja	5
3.1) Uwierzytelnianie	5
3.2) Aktualizacja kolekcji	6
4) Prezentacja działania	7
4.1) Poprawne uruchomienie	7
4.2) Brakujące lub niepoprawne poświadczenia	8
5) Analiza danych	9
6) Wnioski	17

1) Cel projektu

Celem projektu było stworzenie aplikacji importującej dane z dowolnej strony internetowej do baz typu NoSQL. W tym celu należało wykorzystać przynajmniej dwie bazy, tj lokalną oraz udostępnioną w ramach darmowej chmury MongoDB. Dane do minimum jednej z tych baz miały być pozyskiwane na zasadzie web scrapping'u.

2) Zbiór danych

Trzynastu laureatów otrzymało Nagrodę Nobla w 2021 r. za osiągnięcia, które przyniosły ludzkości największe korzyści. Ich praca i odkrycia sięgają od klimatu Ziemi i naszego zmysłu dotyku po wysiłki na rzecz ochrony wolności słowa.

Między innymi te osiągnięcia stanowią zbiór danych użyty w projekcie.

Dane składają się z dwóch kolekcji laureates oraz nobelPrizes, które zostały załadowane do lokalnej bazy danych z ogólnie dostępnego interfejsu API oparty na RES: <https://www.nobelprize.org/about/developer-zone-2/>

Wyniki zapytań do API jest zwracany w formacie JSON, a zawartość jest aktualizowana na bieżąco w oparciu o informację www.nobelprize.org, w momencie ogłoszenia nowych Laureatów.

Łącznie surowe dane liczą 6.5 MB, kolekcja *laureates* ma 968 rekordów, natomiast nobelPrizes 658 rekordów.

Schemat kolekcji prezentuje się następująco:

laureates
> { _id ObjectId(0)
> { birth map(0)
> { death map(0)
✓ { links list(0)
✓ { action String(0)
✓ { class array(0)
✓ { rel String(0)
✓ { title String(0)
✓ { href String(0)
✓ { types String(0)
✓ { nobelPrizes list(0)
> { affiliations list(0)
✓ { awardYear String(0)
> { category map(0)
> { categoryFullName map(0)
✓ { dateAwarded String(0)
> { links list(0)
> { motivation map(0)
✓ { portion String(0)
✓ { prizeAmount Integer
✓ { prizeAmountAdjusted Integer
✓ { prizeStatus String(0)
✓ { sortOrder String(0)
✓ { sameAs array(0)
✓ { wikidata map(0)
✓ { id String(0)
✓ { url String(0)
✓ { wikipedia map(0)
✓ { english String(0)
✓ { slug String(0)
> { familyName map(0)
✓ { fileName String(0)
> { fullName map(0)
✓ { gender String(0)
> { givenName map(0)
✓ { id String(0)
> { knownName map(0)

prizes
✓ { _id ObjectId(0)
✓ { awardYear String(0)
✓ { category map(0)
✓ { en String(0)
✓ { no String(0)
✓ { se String(0)
✓ { laureates list(0)
✓ { id String(0)
✓ { knownName map(0)
✓ { en String(0)
✓ { links list(0)
✓ { action String(0)
✓ { href String(0)
✓ { rel String(0)
✓ { types String(0)
✓ { motivation map(0)
✓ { en String(0)
✓ { no String(0)
✓ { se String(0)
✓ { portion String(0)
✓ { sortOrder String(0)
✓ { links list(0)
✓ { action String(0)
✓ { href String(0)
✓ { rel String(0)
✓ { types String(0)
✓ { prizeAmount Integer
✓ { prizeAmountAdjusted Integer
✓ { categoryFullName map(0)
✓ { en String(0)
✓ { no String(0)
✓ { se String(0)
✓ { dateAwarded String(0)

3) Implementacja

Projekt został zaimplementowany w języku JavaScript.

3.1) Uwierzytelnianie

Obie bazy danych są zabezpieczone hasłem. W programie poświadczenia wymagane do dostępu do baz są wyciągane z pliku konfiguracyjnego. W przypadku ich braku wyświetlany jest komunikat o błędzie.

```
if (process.env.NODE_ENV !== 'production') {  
  require('dotenv').config();  
}  
  
if (!USERNAME || !PASSWORD) {  
  throw new Error('No credentials provided. Please update .env file');  
}  
  
const local_uri = `mongodb://${USERNAME}:${PASSWORD}@localhost:27017`;  
const atlas_uri =  
`mongodb+srv://${USERNAME}:${PASSWORD}@cluster0.mmpd8.mongodb.net/myFirstDatabase?r  
etryWrites=true&w=majority`;
```

Podobnie jeśli poświadczenia okażą się błędne, wyświetlany jest odpowiedni komunikat.

```
} catch(e) {  
  if (e.codeName === 'AuthenticationFailed') {  
    console.error('Authentication failed.');  }  
}
```

3.2) Aktualizacja kolekcji

Po podłączeniu się do obu baz wywoływane są metody aktualizujące dane. Przyjmują one nazwę kolekcji do zaktualizowania, zapytanie do API zwracające aktualne dane, a także połączenie do bazy, w której ma się wykonać aktualizacja.

```
await local_client.connect();
console.log('Connected successfully to local db');
await atlas_client.connect();
console.log('Connected successfully to atlas db');

await updateCollection(laureatesName, axios.get(`${baseApiUrl}${laureatesName}`,
{ params: { limit: 968 } })), local_client);

await updateCollection(laureatesName, axios.get(`${baseApiUrl}${laureatesName}`,
{ params: { limit: 968 } })), atlas_client);
```

Funkcja aktualizująca zaczyna od wykonania przekazanego zapytania. Następnie usuwa wszystkie aktualnie istniejące rekordy z wybranej kolekcji i wypełnia kolekcję nowo otrzymanymi danymi. O każdym z tych kroków informuje użytkownika komunikatami.

```
async function updateCollection(name, call, client) {
  console.log('Started fetching data...');
  const { data } = await call;
  console.log('Finished fetching data.');

  const database = client.db('nobel');
  const collection = database.collection(name);
  await collection.deleteMany();

  console.log('Started inserting records...');
  const result = await collection.insertMany(data[name]);
  console.log(`Finished. Inserted ${data[name].length} records to collection
${name}.`);
}
```

Po zaktualizowaniu danych pobieranych z API następuje aktualizacja tabeli statystyk. Wykonywana jest seria zapytań do bazy lokalnej a ich wyniki są umieszczane w odpowiedniej kolekcji bazy chmurowej.

```
await updateStatsCollection(local_client, atlas_client);
```

```

async function updateStatsCollection(sourceClient, destinationClient) {
  console.log('Started filling statistics collection...');
  const source_database = sourceClient.db('nobel');
  const laureates = source_database.collection(laureatesName);
  const numberOfMen = await laureates.find({ gender: 'male' }).count();
  const numberOfWomen = await laureates.find({ gender: 'female' }).count();
  const numberOfPoles = await laureates.find({ "birth.place.country.en": "Poland" }).count();
  const avgAgeFirstPrize = await laureates.aggregate([...
]).toArray();
  const numberOfMultiWinners = await laureates.aggregate([...
]).toArray();

  const destination_database = destinationClient.db('nobel');
  const stats = destination_database.collection(statsName);
  await stats.deleteMany();
  await stats.insertOne({ name: 'Male count', value: numberOfMen });
  await stats.insertOne({ name: 'Female count', value: numberOfWomen });
  await stats.insertOne({ name: 'Pole count', value: numberOfPoles });
  await stats.insertOne({ name: 'Average age when getting first prize', value: avgAgeFirstPrize[0].avgAge });
  await stats.insertOne({ name: 'Multiple awarded count', value: numberOfMultiWinners[0].count });
  console.log('Finished filling statistics collection.');
```

4) Prezentacja działania

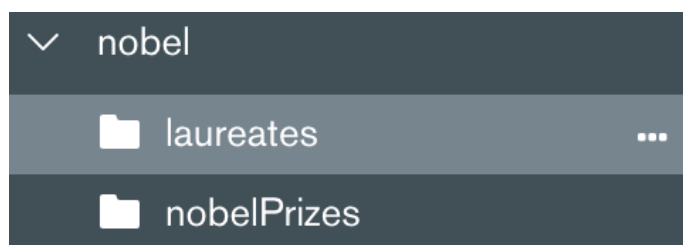
4.1) Poprawne uruchomienie

W przypadku poprawnego uruchomienia wyświetlane są komunikaty o kolejnych krokach działania.

```

gcks@gcks projekt 1 % node index.js
Connected successfully to local db
Connected successfully to atlas db
Started fetching data...
Finished fetching data.
Started inserting records...
Finished. Inserted 968 records to collection laureates.
Started fetching data...
Finished fetching data.
Started inserting records...
Finished. Inserted 658 records to collection nobelPrizes.
Started fetching data...
Finished fetching data.
Started inserting records...
Finished. Inserted 968 records to collection laureates.
Started fetching data...
Finished fetching data.
Started inserting records...
Finished. Inserted 658 records to collection nobelPrizes.
Started filling statistics collection...
Finished filling statistics collection.
Closed connection to local db
Closed connection to atlas db
```

Dane w obu bazach są aktualizowane.



FILTER { field: 'value' } **OPTIONS** **FIND** **RESET** **REFRESH**

ADD DATA **VIEW** **REFRESH**

Displaying documents 1 - 20 of 968

```
{
  "_id": ObjectId("6197e847f197dc00ca3f948a"),
  "id": "745",
  "knownName": Object,
  "givenName": Object,
  "familyName": Object,
  "fullName": Object,
  "fileName": "spence",
  "gender": "male",
  "birth": Object,
  "wikipedia": Object,
  "wikidata": Object,
  "sameAs": Array,
  "links": Array,
  "nobelPrizes": Array
}
```

+ Create Database

Q NAMESPACES

▼ nobel

laureates

nobelPrizes

stats

nobel.stats

COLLECTION SIZE: 314B TOTAL DOCUMENTS: 5 INDEXES TOTAL SIZE: 36KB

Find

Indexes

Schema Anti-Patterns 0

Aggregation

Search Indexes ●

FILTER { field: 'value' } **OPTIONS**

QUERY RESULTS 1-5 OF 5

```
{
  "_id": ObjectId("61a3dbfe619ae2c9a06d8845"),
  "name": "Male count",
  "value": 885
}
```

```
{
  "_id": ObjectId("61a3dbfe619ae2c9a06d8846"),
  "name": "Female count",
  "value": 58
}
```

4.2) Brakujące lub niepoprawne poświadczenia

W przypadku nieudanego uwierzytelniania aplikacja również wyświetla stosowne informacje o błędach.

```
gcks@gcks projekt 1 % node index.js
No credentials provided. Please update .env file
Authentication failed.
```

```
gcks@gcks projekt 1 % node index.js
Authentication failed.
```


5) Analiza danych

W ramach korzystania z bazy danych MongoDB, możemy posługiwać się dedykowanymi operacjami filtrowania danych dla własnej analizy. Między innymi dzięki użyciu agregacji możemy przetwarzać dokumenty i zwracać przeliczone wyniki.

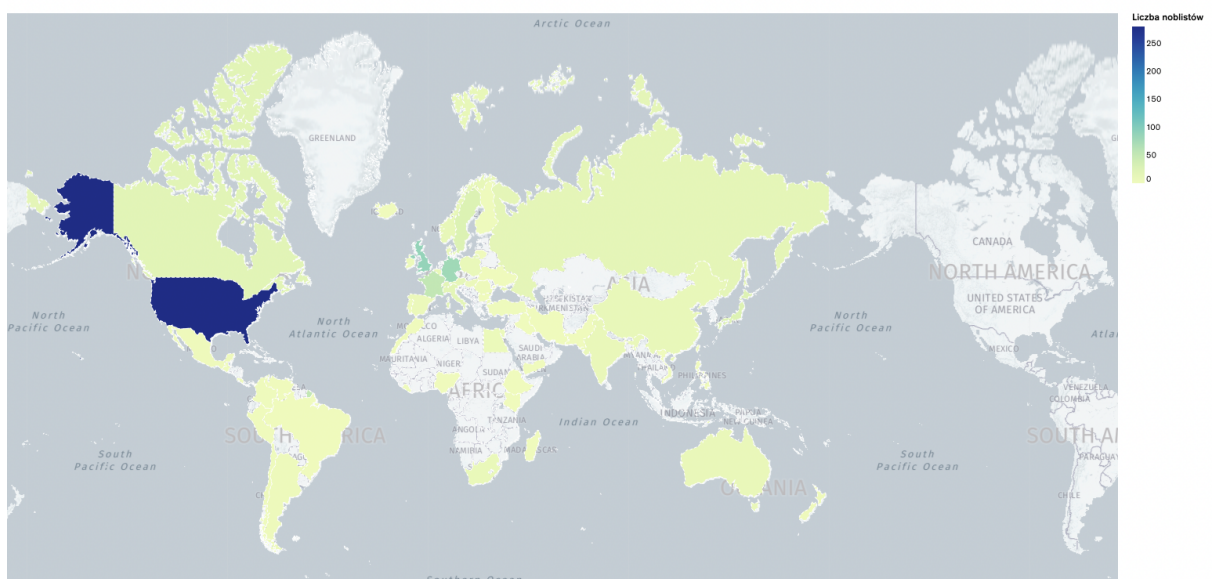
Poniżej zostały zaprezentowane przykładowe statystyki, przy wykorzystaniu technik jakie udostępnia baza MongoDB.

- 1) Liczba nagród z podziałem na kraje pochodzenia laureatów.

```
Laureates.aggregate( pipeline: [  
  {$unwind: "$nobelPrizes"},  
  {  
    $group: {  
      _id: "$birth.place.country.en",  
      sum: {$sum: 1},  
    },  
  },  
  {  
    $sort: {sum: -1}  
  }  
])
```

```
1 [   
2 {   
3   "_id": "USA",   
4   "sum": 283   
5 },   
6 {   
7   "_id": "United Kingdom",   
8   "sum": 89   
9 },   
10 {   
11   "_id": "Germany",   
12   "sum": 80   
13 },   
14 {   
15   "_id": "France",   
16   "sum": 54   
17 },   
18 {   
19   "_id": null,   
20   "sum": 30   
21 },   
22 {   
23   "_id": "Sweden",   
24   "sum": 29   
25 },   
26 {   
27   "_id": "Japan",   
28   "sum": 28   
29 },   
30 {   
31   "_id": "Canada",   
32   "sum": 21   
33 },   
34 {   
35   "_id": "the Netherlands",   
36   "sum": 19   
37 }
```

Liczba noblistów wg miejsca urodzenia

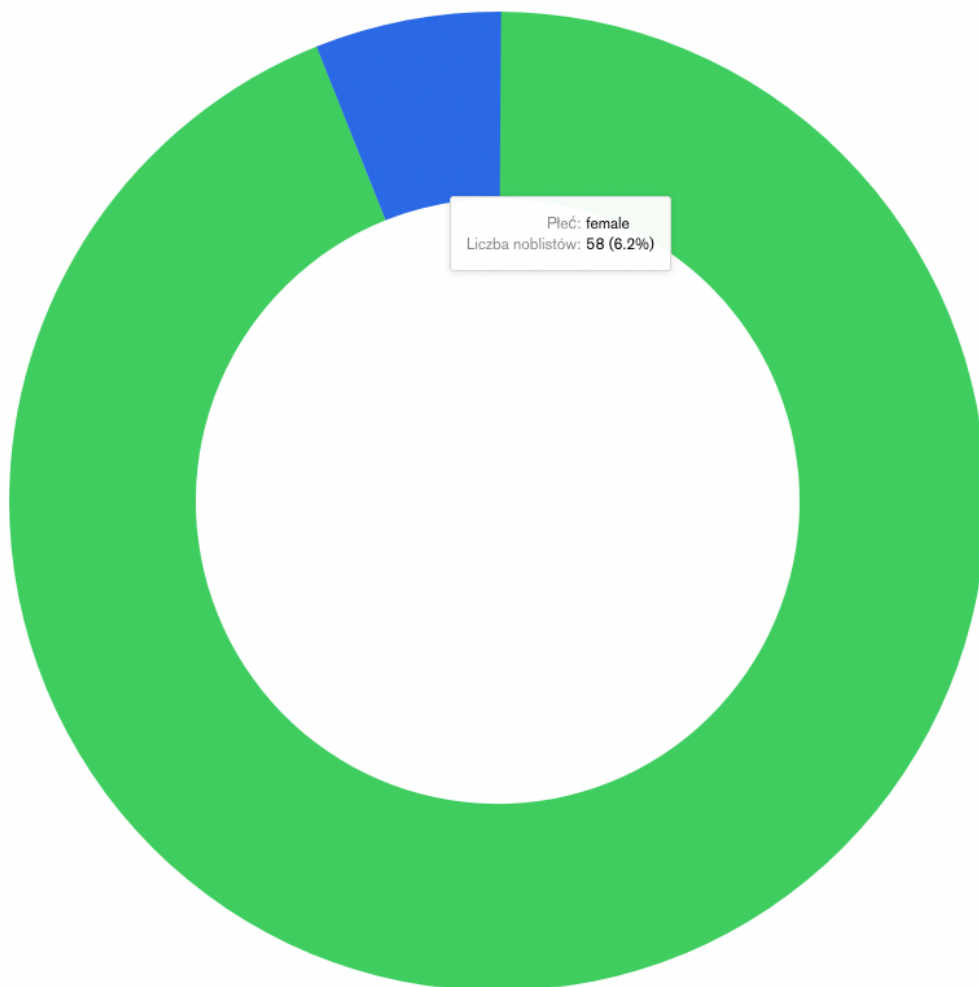


2) Liczba nagród z podziałem na płeć laureatów.

```
Laureates.aggregate( pipeline: [
  {$unwind: "$nobelPrizes"},
  {
    $group: {
      _id: "$gender",
      sum: {$sum: 1},
    }
  },
  {
    $sort: {sum: -1}
  },
  {
    $limit: 100
  }
])
```

```
1 [
2   {
3     "_id": "male",
4     "sum": 888
5   },
6   {
7     "_id": "female",
8     "sum": 59
9   },
10  {
11    "_id": null,
12    "sum": 28
13  }
14 ]
```

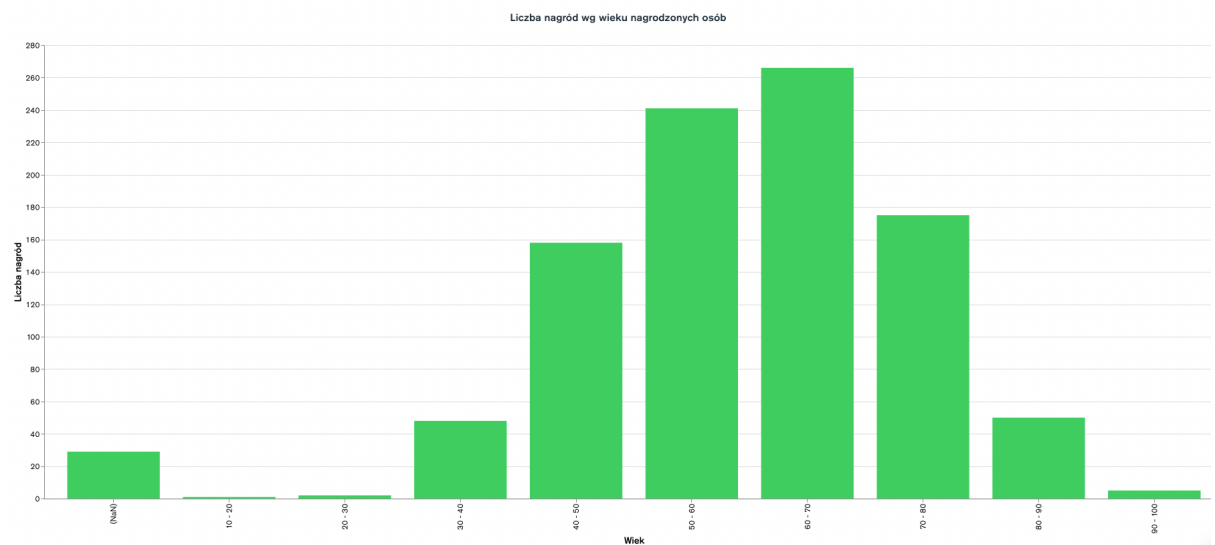
Liczba noblistów wg płci



3) Liczba nagród z podziałem na wiek laureatów.

```
Laureates.aggregate( pipeline: [
  {$unwind: "$nobelPrizes"},
  {
    $project: {
      _id: 1,
      "birthYear": {
        $convert: {
          input: {$substr: ["$birth.date", 0, 4 ]},
          to: "int",
          onNull: 0
        } },
      "awardYear": {
        $convert: {
          input: "$nobelPrizes.awardYear", to: "int",
          onNull: 0
        } },
      "gender": 1
    }
  },
  {
    $addFields: {
      age: {$subtract: [ "$awardYear", "$birthYear" ]},
    }
  },
  {
    $group: {
      _id: '$age',
      count: {$sum: 1},
    }
  },
  {
    $sort: {count: -1}
  }
])
```

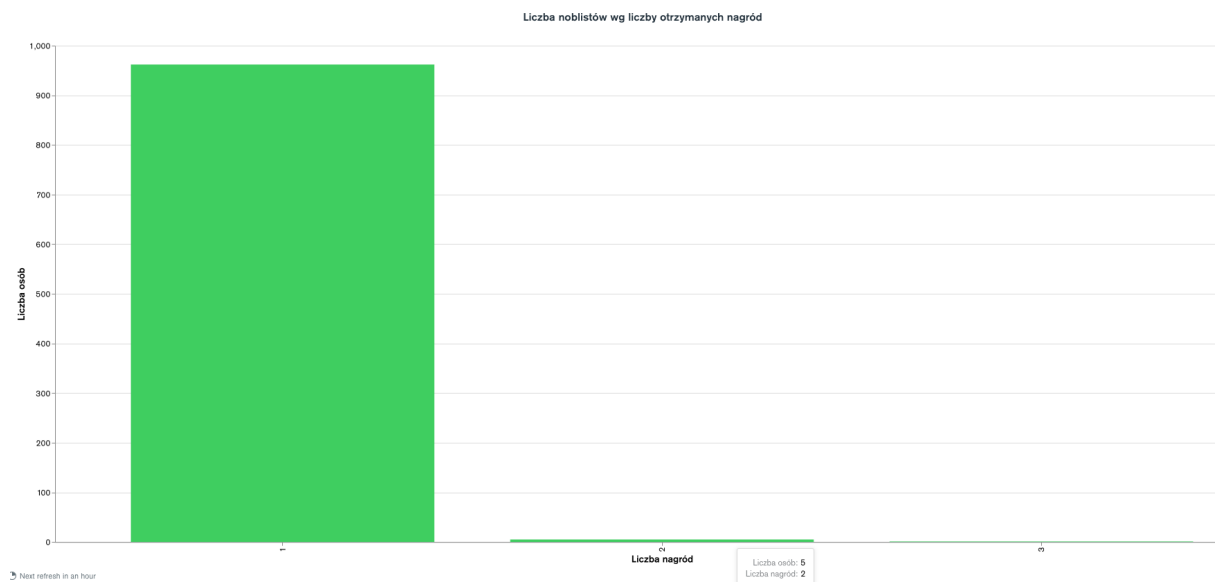
```
1 [
2 {
3   "_id": 63,
4   "count": 35
5 },
6 {
7   "_id": 61,
8   "count": 34
9 },
10 {
11   "_id": 56,
12   "count": 34
13 },
14 {
15   "_id": 60,
16   "count": 32
17 },
18 {
19   "_id": 64,
20   "count": 30
21 },
22 {
23   "_id": 54,
24   "count": 30
25 },
26 {
27   "_id": 68,
28   "count": 30
29 },
30 {
31   "_id": null,
32   "count": 29
33 },
34 {
35   "_id": 62,
36   "count": 27
37 },
38 {
39   "_id": 55,
40   "count": 25
41 },
42 {
43   "_id": 57,
44   "count": 25
45 }
```



4) Liczba laureatów z podziałem na liczbę nagród.

```
Laureates.aggregate( pipeline: [
  {
    $project: {
      _id: 1,
      countOfPrize: {$size: "$nobelPrizes"}
    }
  },
  {
    $group: {
      _id: "$countOfPrize",
      count : { $sum: 1}
    }
  }
])
```

```
1  [
2  {
3    "id": 2,
4    "count": 5
5  },
6  {
7    "id": 3,
8    "count": 1
9  },
10 {
11    "id": 1,
12    "count": 962
13  }
14 ]
```



5) Średnia wieku laureatów dla ich pierwszego zdobytego tytułu.

```
Laureates.aggregate( pipeline: [
  {
    $project: {
      _id: 1,
      "birthYear": { $convert: {
        input: { $substr: ["$birth.date", 0, 4]},
        to: "int",
        onNull: 0
      } },
      "awardYear": { $convert: {
        input: { $first: "$nobelPrizes.awardYear"},
        to: "int",
        onNull: 0
      } },
    }
  },
  {
    $addFields: {
      age: { $subtract: [ "$awardYear", "$birthYear" ] },
    }
  },
  {
    $group: {
      _id: '_id',
      avgAge: { $avg: "$age" },
    }
  },
],
})
```

```
1 [
2 {
3   "_id": "_id",
4   "avgAge": 60.03503184713376
5 }
6 ]
```

6) Liczba laureatów z wieloma nagrodami.

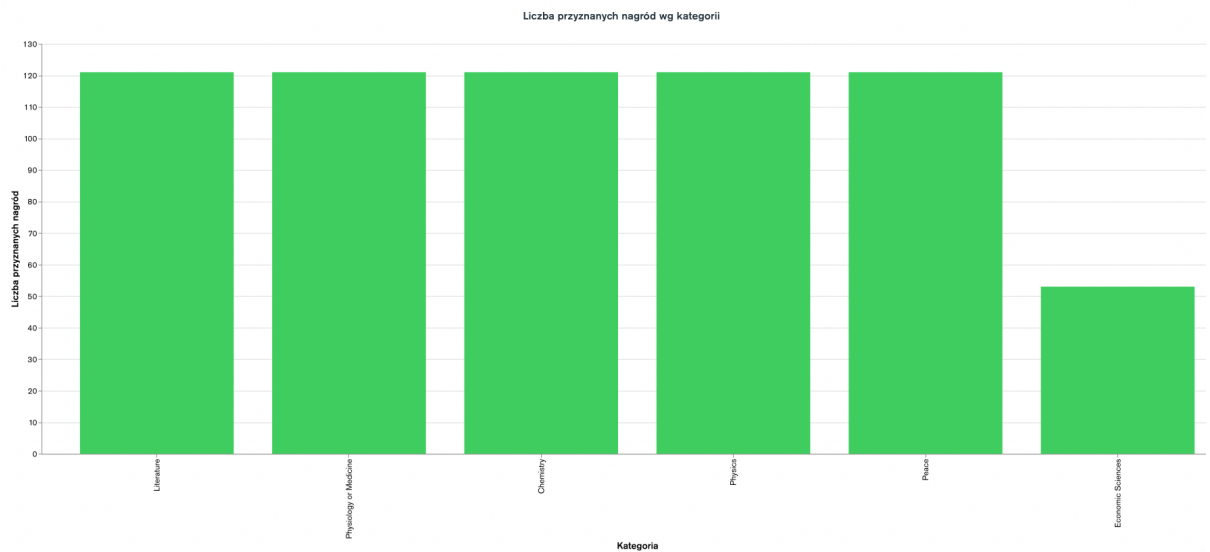
```
Laureates.aggregate( pipeline: [
  {
    $project: {
      _id: 1,
      countOfPrize: { $size: "$nobelPrizes" }
    }
  },
  {
    $match: {
      countOfPrize : { $gt: 1 }
    }
  },
  {
    $group: {
      _id: "_id",
      count : { $sum: 1 }
    }
  },
],
})
```

```
1 [
2 {
3   "_id": "_id",
4   "count": 6
5 }
6 ]
```

7) Nagrody Nobla z podziałem na kategorie.

```
NobelPrizes.aggregate( pipeline: [
  {
    $group: {
      _id: "$category.en",
      sum: {$sum: 1},
    }
  },
  {
    $sort: {sum: -1}
  }
])
```

```
1 [
2   {
3     "_id": "Peace",
4     "sum": 121
5   },
6   {
7     "_id": "Physiology or Medicine",
8     "sum": 121
9   },
10  {
11    "_id": "Chemistry",
12    "sum": 121
13  },
14  {
15    "_id": "Literature",
16    "sum": 121
17  },
18  {
19    "_id": "Physics",
20    "sum": 121
21  },
22  {
23    "_id": "Economic Sciences",
24    "sum": 53
25  }
26 ]
```



8) Średnia przyznawana kwota finansowa z podziałem na kategorie.

```
NobelPrizes.aggregate( pipeline: [  
  {  
    $group: {  
      _id: "$category.en",  
      avgOfPrizeAmount:  
        {$avg: "$prizeAmount"},  
      avgOfPrizeAmountAdjusted:  
        {$avg: "$prizeAmountAdjusted"}  
    }  
  },  
  {  
    $sort: {avgOfPrizeAmount: -1}  
  }  
])
```

```
1  [  
2  {  
3    "_id": "Economic Sciences",  
4    "avgOfPrizeAmount": 5618037.735849056,  
5    "avgOfPrizeAmountAdjusted": 7490924.886792453  
6  },  
7  {  
8    "_id": "Physics",  
9    "avgOfPrizeAmount": 2554777.5454545454,  
10   "avgOfPrizeAmountAdjusted": 5702950.289256198  
11  },  
12  {  
13   "_id": "Literature",  
14   "avgOfPrizeAmount": 2554777.5454545454,  
15   "avgOfPrizeAmountAdjusted": 5701601.52892562  
16  },  
17  {  
18   "_id": "Physiology or Medicine",  
19   "avgOfPrizeAmount": 2554777.5454545454,  
20   "avgOfPrizeAmountAdjusted": 5702950.289256198  
21  },  
22  {  
23   "_id": "Chemistry",  
24   "avgOfPrizeAmount": 2554777.5454545454,  
25   "avgOfPrizeAmountAdjusted": 5702950.289256198  
26  },  
27  {  
28   "_id": "Peace",  
29   "avgOfPrizeAmount": 2554777.5454545454,  
30   "avgOfPrizeAmountAdjusted": 5702950.289256198  
31  }  
32 ]
```

9) Liczba przyznanych Nagród Nobla z podziałem na lata.

```
NobelPrizes.aggregate( pipeline: [  
  {  
    $project: {  
      _id: 1,  
      "dateAwarded": {  
        $convert: {  
          input: {$substr:  
            ["$dateAwarded", 0, 4]  
          },  
          to: "int",  
          onNull: 0  
        }  
      }  
    },  
  },  
  {  
    $group: {  
      _id: "$dateAwarded",  
      sum: {$sum: 1},  
    },  
  },  
  {  
    $sort: {_id: -1}  
  }  
])
```

```
1 [  
2 {  
3   "_id": 2021,  
4   "sum": 6  
5 },  
6 {  
7   "_id": 2020,  
8   "sum": 6  
9 },  
10 {  
11   "_id": 2019,  
12   "sum": 7  
13 },  
14 {  
15   "_id": 2018,  
16   "sum": 3  
17 },  
18 {  
19   "_id": 2017,  
20   "sum": 5  
21 },  
22 {  
23   "_id": 2016,  
24   "sum": 5  
25 },  
26 {  
27   "_id": 2015,  
28   "sum": 5  
29 },  
30 {  
31   "_id": 2014,  
32   "sum": 5  
33 },  
34 {  
35   "_id": 2013,  
36   "sum": 5  
37 },  
38 {  
39   "_id": 2012,  
40   "sum": 5  
41 },  
42 ]
```


6) Wnioski

Zakres funkcjonalności i zastosowanych technik został wykonany zgodnie z założeniami wymaganymi na ocenę 4.0

Z analizy danych wynika, że najczęściej laureatów nagrody Nobla pochodzi z krajów rozwiniętych o znacznej liczbie ludności takich jak Stany Zjednoczone, Wielka Brytania czy Niemcy.

Nagroda Nobla znacznie częściej otrzymywali mężczyźni niż kobiety oraz osoby zaawansowane wiekowo. Najwięcej nagród zostało przyznanych w grupie 60-70 lat. Średni wiek laureata nagrody Nobla to około 60 lat.

Niewielu laureatom nagrody udało się ponownie ją zdobyć, udało się to tylko 6 osobom.

Kwoty przyznawane w ramach nagrody Nobla, były największe w przypadku nauk ekonomicznych, w pozostałych kategoriach ich wartość nie różniła się istotnie. Należy jednak zauważyć, że liczba przyznanych nagród w przypadku nauk ekonomicznych była najmniejsza spośród wszystkich kategorii.