



Zanim zaczniemy ćwiczenia musimy zainstalować:

- PIP:  
`yum install python-pip`
- MRJob:  
`pip install google-api-python-client==1.6.4`  
`pip install mrjob==0.5.11`
- Nano:  
`yum install nano`

oraz ściągnąć nasz zbiór danych:

- Utwórz nowy folder w lokalnym systemie plików:  
`mkdir twoje_nazwisko`
- Wejdź do folderu:  
`cd twoje_nazwisko`
- Pobierz plik:  
`wget https://danilewicz.blob.core.windows.net/public/1987.zip`
- Rozpakuj archiwum:  
`unzip 1987.zip`

	Name	Description
1	Year	1987-2008
2	Month	01-Dec
3	DayofMonth	Jan-31
4	DayOfWeek	1 (Monday) - 7 (Sunday)
5	DepTime	actual departure time (local, hhmm)
6	CRSDepTime	scheduled departure time (local, hhmm)
7	ArrTime	actual arrival time (local, hhmm)
8	CRSArrTime	scheduled arrival time (local, hhmm)
9	UniqueCarrier	unique carrier code
10	FlightNum	flight number
11	TailNum	plane tail number
12	ActualElapsedTime	in minutes
13	CRSElapsedTime	in minutes
14	AirTime	in minutes
15	ArrDelay	arrival delay, in minutes
16	DepDelay	departure delay, in minutes
17	Origin	origin IATA airport code
18	Dest	destination IATA airport code
19	Distance	in miles
20	TaxiIn	taxi in time, in minutes

21	TaxiOut	taxi out time in minutes
22	Cancelled	was the flight cancelled?
23	CancellationCode	reason for cancellation (A = carrier, B = weather, C = NAS, D = security)
24	Diverted	1 = yes, 0 = no
25	CarrierDelay	in minutes
26	WeatherDelay	in minutes
27	NASDelay	in minutes
28	SecurityDelay	in minutes
29	LateAircraftDelay	in minutes

## Ćwiczenie 1

Stworzymy prosty program MapReduce, który zliczy liczbę lotów w poszczególnych miesiącach.

1. Stwórz nowy plik w którym będzie nasz kod:  
`nano liczbilotow.py`
2. W edytorze tekstu wpisz poniższy kod:

```
from mrjob.job import MRJob
from mrjob.step import MRStep

class LiczbaLotow(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper_get_months,
                    reducer=self.reducer_count_months)
        ]

    def mapper_get_months(self, _, line):
        fields = line.split(',')
        yield fields[1], 1

    def reducer_count_months(self, key, values):
        yield key, sum(values)

if __name__ == '__main__':
    LiczbaLotow.run()
```

3. Zamknij edytor zapisując zmiany:
  - `Ctrl + X`
  - `Y`
  - `Enter`

4. Najpierw uruchomimy nasz program lokalnie: `python liczbilotow.py 1987.csv`

```
maria_dev@sandbox-hdp:~/danilewicz
[maria_dev@sandbox-hdp danilewicz]$ nano liczbilotow.py
[maria_dev@sandbox-hdp danilewicz]$ python liczbilotow.py 1987.csv
No configs found; falling back on auto-configuration
Creating temp directory /tmp/liczbilotow.maria_dev.20190208.090447.060327
Running step 1 of 1...
Streaming final output from /tmp/liczbilotow.maria_dev.20190208.090447.060327/output...
"12"      440403
"Month"   1
"10"      448620
"11"      422803
Removing temp directory /tmp/liczbilotow.maria_dev.20190208.090447.060327...
[maria_dev@sandbox-hdp danilewicz]$
```

5. Następnie uruchomimy nasz program korzystając z hadoop:  
`python liczbilotow.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar 1987.csv`

```
maria_dev@sandbox-hdp:~/danilewicz
Reduce input records=1311827
Reduce output records=4
Reduce shuffle bytes=11806458
Shuffled Maps =2
Spilled Records=2623654
Total committed heap usage (bytes)=284164096
Virtual memory (bytes) snapshot=6420058112
Shuffle Errors
BAD_ID=0
CONNECTION=0
IO_ERROR=0
WRONG_LENGTH=0
WRONG_MAP=0
WRONG_REDUCE=0
Streaming final output from hdfs:///user/maria_dev/tmp/mrjob/liczbilotow.maria_dev.20190208.131501.424085/output...
"10"      448620
"11"      422803
"12"      440403
"Month"   1
Removing HDFS temp directory hdfs:///user/maria_dev/tmp/mrjob/liczbilotow.maria_dev.20190208.131501.424085...
Removing temp directory /tmp/liczbilotow.maria_dev.20190208.131501.424085...
[maria_dev@sandbox-hdp danilewicz]$
```

## Ćwiczenie 2

Stworzymy program, który znajdzie nam lotniska z których startowało najwięcej samolotów.

1. Zaczniemy od stworzenia programu podobnego do ćwiczenia 1:  
Stwórz nowy plik w którym będzie nasz kod:  
[nano lotniska.py](#)
2. W edytorze tekstu wpisz poniższy kod:

```
from mrjob.job import MRJob
from mrjob.step import MRStep

class Lotniska(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper_get_airports,
                    reducer=self.reducer_count_airports)
        ]

    def mapper_get_airports(self, _, line):
        fields = line.split(',')
        yield fields[16], 1

    def reducer_count_airports(self, key, values):
        yield key, sum(values)

if __name__ == '__main__':
    Lotniska.run()
```

3. Zamknij edytor zapisując zmiany:
  - *Ctrl + X*
  - *Y*
  - *Enter*

4. Uruchom program lokalnie: `python lotniska.py 1987.csv`

```
maria_dev@sandbox-hdp:~/danilewicz
"HSV" 1546
"HTS" 620
"IAD" 14560
"IAH" 21566
"ICT" 3575
"IDA" 326
"ILG" 29
"ILM" 1150
"IND" 8817
"ISO" 367
"ISP" 2050
"ITH" 640
"JAC" 226
"JAN" 2583
"JAX" 6498
"JFK" 12273
"JNU" 728
"KOA" 159
"KTN" 508
"LAN" 734
"LAS" 19239
"LAX" 45646
Removing temp directory /tmp/lotniska.maria_dev.20190208.204919.712077...
[maria_dev@sandbox-hdp danilewicz]$
```

5. Zamień kolejność w `reducer_count_airports`:

```
from mrjob.job import MRJob
from mrjob.step import MRStep

class Lotniska(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper_get_airports,
                    reducer=self.reducer_count_airports)
        ]

    def mapper_get_airports(self, _, line):
        fields = line.split(',')
        yield fields[16], 1

    def reducer_count_airports(self, key, values):
        yield sum(values), key

if __name__ == '__main__':
    Lotniska.run()
```

6. Sprawdźmy teraz wynik: [python lotniska.py 1987.csv](#)

```
maria_dev@sandbox-hdp:~/danilewicz
1546      "HSV"
620       "HTS"
14560     "IAD"
21566     "IAH"
3575      "ICT"
326       "IDA"
29        "ILG"
1150      "ILM"
8817      "IND"
367       "ISO"
2050      "ISP"
640       "ITH"
226       "JAC"
2583      "JAN"
6498      "JAX"
12273     "JFK"
728       "JNU"
159       "KOA"
508       "KTN"
734       "LAN"
19239     "LAS"
45646     "LAX"

Removing temp directory /tmp/lotniska.maria_dev.20190208.205759.251583...
[maria_dev@sandbox-hdp danilewicz]$
```

7. Aby ponownie posortować wynik, dopisz liniki zaznaczone za żółto:

```
from mrjob.job import MRJob
from mrjob.step import MRStep

class Lotniska(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper_get_airports,
                   reducer=self.reducer_count_airports),
            MRStep(reducer=self.reducer_sort_output)
        ]

    def mapper_get_airports(self, _, line):
        fields = line.split(',')
        yield fields[16], 1

    def reducer_count_airports(self, key, values):
        yield sum(values), key

    def reducer_sort_output(self, count, airports):
        for airport in airports:
            yield airport, count

if __name__ == '__main__':
    Lotniska.run()
```

8. Sprawdźmy teraz wynik: [python lotniska.py 1987.csv](#)

```
maria_dev@sandbox-hdp:~/danilewicz
"BOS" 25250
"JAN" 2583
"SRQ" 2598
"APF" 261
"HLN" 268
"PHF" 268
"FNT" 274
"SPN" 275
"DTW" 27548
"SIT" 282
"GEG" 2848
"LGA" 28596
"PIT" 28765
"GUC" 29
"ILG" 29
"ALB" 2922
"SAV" 2984
"PHX" 29848
"LFT" 299
"PVD" 3043
"CAE" 3058
"EWB" 30991
Removing temp directory /tmp/lotniska.maria_dev.20190208.211805.409695...
[maria_dev@sandbox-hdp danilewicz]$
```

9. Wynik jest źle posortowany bo liczby są posortowane jak tekst. Naprawimy to zmieniając jedną linijkę:

```
from mrjob.job import MRJob
from mrjob.step import MRStep

class Lotniska(MRJob):
    def steps(self):
        return [
            MRStep(mapper=self.mapper_get_airports,
                    reducer=self.reducer_count_airports),
            MRStep(reducer=self.reducer_sort_output)
        ]

    def mapper_get_airports(self, _, line):
        fields = line.split(',')
        yield fields[16], 1

    def reducer_count_airports(self, key, values):
        yield str(sum(values)).zfill(5), key

    def reducer_sort_output(self, count, airports):
        for airport in airports:
            yield airport, count

if __name__ == '__main__':
    Lotniska.run()
```



10. Teraz możemy uruchomić nasz program na Hadoop:

```
python lotniska.py -r hadoop --hadoop-streaming-jar /usr/hdp/current/hadoop-mapreduce-client/hadoop-streaming.jar 1987.csv
```

```
maria_dev@sandbox-hdp:~/danilewicz
"MIA"      "17649"
"MEM"      "19081"
"LAS"      "19239"
"PHL"      "20570"
"IAH"      "21566"
"DCA"      "22016"
"MSP"      "23108"
"CLT"      "24518"
"BOS"      "25250"
"DTW"      "27548"
"LGA"      "28596"
"PIT"      "28765"
"PHX"      "29848"
"EWB"      "30991"
"STL"      "32097"
"SFO"      "35155"
"DEN"      "43376"
"LAX"      "45646"
"DFW"      "51860"
"ATL"      "66309"
"ORD"      "67216"
Removing HDFS temp directory hdfs:///user/maria_dev/tmp/mrjob/lotniska.maria_dev
Removing temp directory /tmp/lotniska.maria_dev.20190208.215925.013731...
[maria_dev@sandbox-hdp danilewicz]$
```

Najwięcej samolotów wystartowało z lotniska ORD – Chicago O'Hare International Airport.