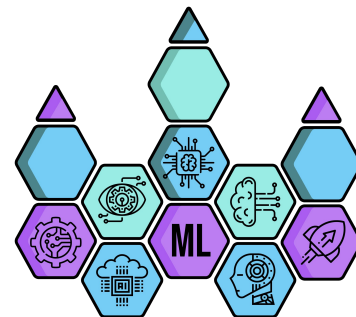# DW Poznań - Duże zbiory danych #9

2019-05-21
Czyli jak ogarnąć dane New York City Crimes

# Agenda

# New York City crimes

- **Dane zawierają informacje na temat przestępczości w Nowym Jorku (New York City Police Department (NYPD) ) za okres  2006 do końca 2017**

- **Mają 7,309,655 rekordów**

- **Plik .csv ma ponad 2 GB**

# Odczyt performance %timeit, %time

```
1 %timeit sum(range(100))
```
```
100000 loops, best of 3: 1.54 µs per loop
```

```
1 %%timeit -n 15
2 data['OFNS_DESC'].value_counts()
```
```
15 loops, best of 3: 783 ms per loop
```

Uruchomienie kilka razy i wyliczenie
średniej czasu uruchomienia

```
1 %%time
2 data.to_csv('rows_2.csv')
```
```
CPU times: user 2min 13s, sys: 3.69 s, total: 2min 16s
Wall time: 2min 17s
```

Wyliczenie czasu uruchomienia
komórki

# %PRUN

Wyliczenie czasu uruchomienia
pojedynczych wierszy

```
1 def sum_of_lists(N):
2     total = 0
3     for i in range(5):
4         L = [j ^ (j >> i) for j in range(N)]
5         total += sum(L)
6     return total
```

```
1 %prun sum_of_lists(1000000)
```

```
        14 function calls in 0.674 seconds

   Ordered by: internal time

   ncalls  tottime  percall  cumtime  percall filename:lineno(function)
        5    0.607    0.121    0.607    0.121 <ipython-input-29-f105717832a2>:4(<listcomp>)
        1    0.030    0.030    0.665    0.665 <ipython-input-29-f105717832a2>:1(sum_of_lists)
        5    0.028    0.006    0.028    0.006 {built-in method builtins.sum}
        1    0.009    0.009    0.674    0.674 <string>:1(<module>)
        1    0.000    0.000    0.674    0.674 {built-in method builtins.exec}
        1    0.000    0.000    0.000    0.000 {method 'disable' of '_lsprof.Profiler' objects}
```
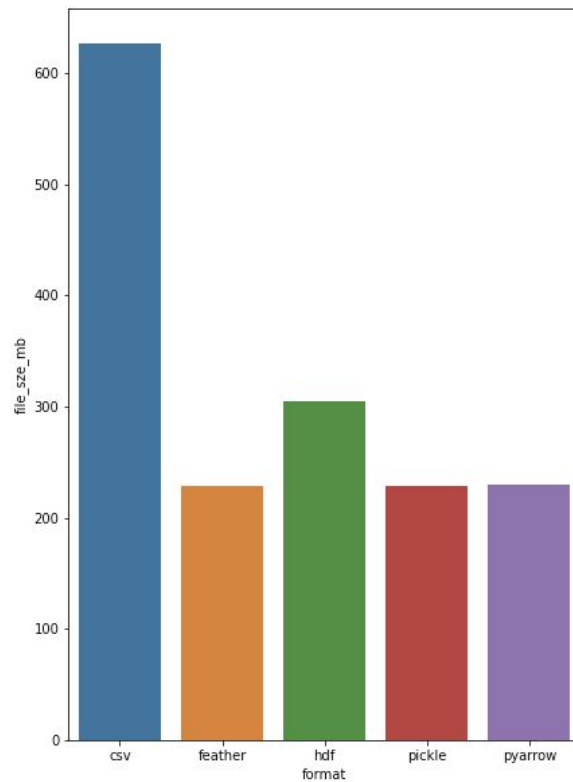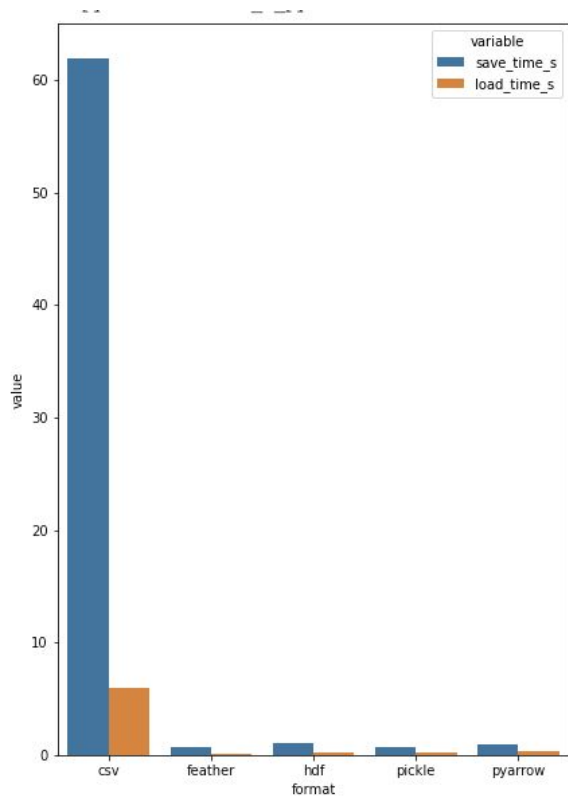
# Memory_profiler %memit

Wyliczenie pamięć użytą podczas operacji

```python
1  def sum_of_lists(N):
2      total = 0
3      for i in range(5):
4          L = [j ^ (j >> i) for j in range(N)]
5          total += sum(L)
6      return total
```
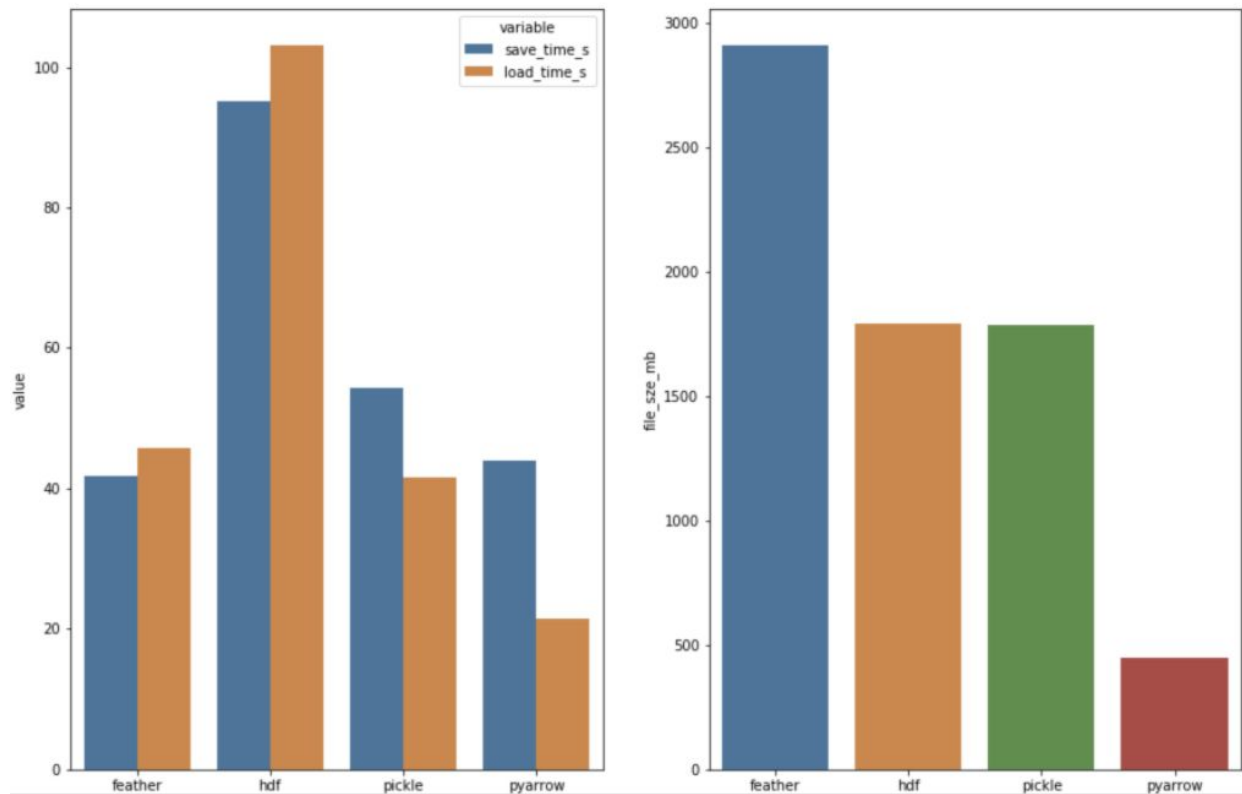
```python
1  %prun sum_of_lists(1000000)
```

# Csv - DataFrame: 10000000 rows

https://towardsdatascience.com/the-best-format-to-save-pandas-data-414dca023e0d

# Pickle, hdf5, feathrow, PyArrow - NYC Data

# Dask

*Dask is a flexible library for parallel computing in Python.*

Dask is composed of two parts:

1. **Dynamic task scheduling** optimized for computation. This is similar to *Airflow, Luigi, Celery, or Make*, but optimized for interactive computational workloads.
2. **"Big Data" collections** like parallel arrays, dataframes, and lists that extend common interfaces like *NumPy, Pandas, or Python iterators* to larger-than-memory or distributed environments. These parallel collections run on top of dynamic task schedulers.

https://docs.dask.org/en/latest/

```python
import pandas as pd
df = pd.read_csv('2015-01-01.csv')
df.groupby(df.user_id).value.mean()
```

```python
import dask.dataframe as dd
df = dd.read_csv('2015-*-*.csv')
df.groupby(df.user_id).value.mean().compute()
```
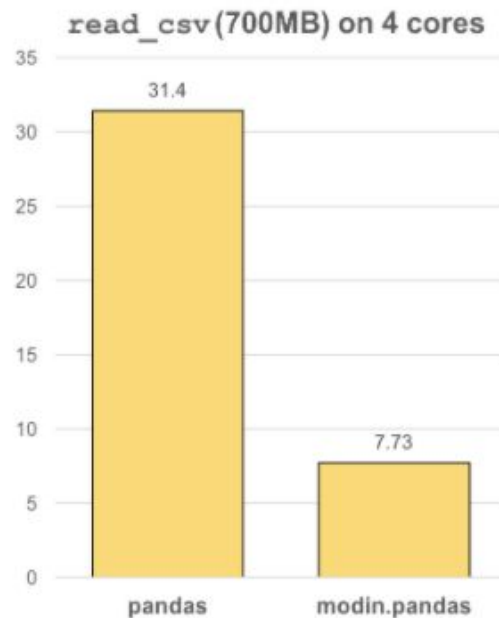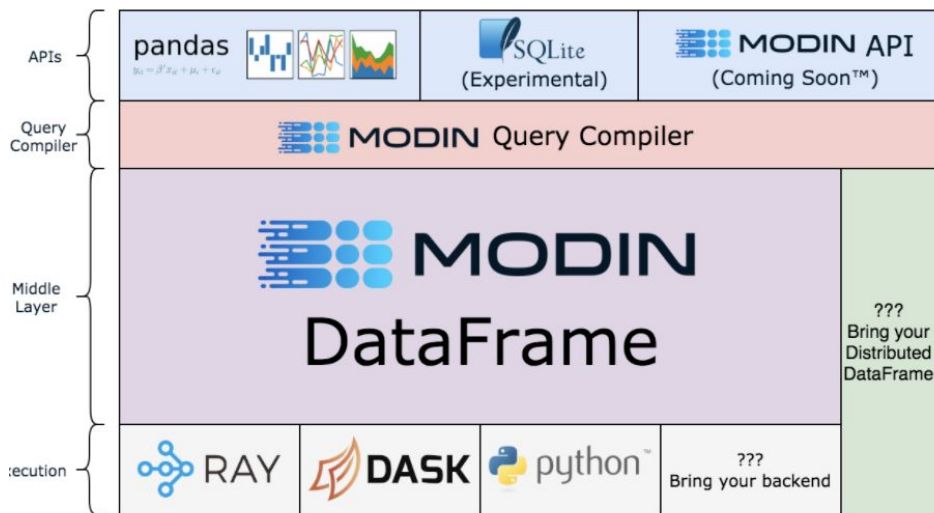
# MODIN

- Jest to multiprocesorowy Dataframe który może przyspieszyć wiele operacji.

# CUDF

- Biblioteka która cały DataFrame wysyła do pamięci GPU gdzie jest przetwarzany.
- Wszystko przez parsowanie CSV do parsowania jest robione po stronie GPU

# Kolejne kroki

- Spotkanie za 2 tygodnie
- Projekt

https://github.com/dataworkshop/dw-poznan-project

# Dziękuję