# Sprawozdanie 10

**Grupa A3:**

inż. Michał Liss

inż. Marceli Sokólski

inż. Piotr Krzystanek

## Instalacja spark w kontenerze namenode

```
RUN tar -xzf spark-3.4.0-bin-hadoop3-scala2.13.tgz
RUN mv spark-3.4.0-bin-hadoop3-scala2.13 /usr/local/spark
RUN rm spark-3.4.0-bin-hadoop3-scala2.13.tgz
ENV SPARK_HOME=/usr/local/spark
ENV PATH="${PATH}:${SPARK_HOME}/bin"
```

## Instalacja python w pozostałych węzłach klastra

```
RUN apt-get update && DEBIAN_FRONTEND=noninteractive apt-get install -y --
no-install-recommends python3.6
RUN update-alternatives --install /usr/bin/python python /usr/bin/python3 1
```

## pyspark

```
○ ❯ docker exec -it namenode bash
root@96cd8ec924c2:~# pyspark
Python 3.7.3 (default, Oct 31 2022, 14:04:00)
[GCC 8.3.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
23/06/12 10:14:49 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
/usr/local/spark/python/pyspark/context.py:317: FutureWarning: Python 3.7 support is deprecated in Spark 3.4.
  warnings.warn("Python 3.7 support is deprecated in Spark 3.4.", FutureWarning)
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 3.4.0
      /_/

Using Python version 3.7.3 (default, Oct 31 2022 14:04:00)
Spark context Web UI available at http://96cd8ec924c2:4040
Spark context available as 'sc' (master = local[*], app id = local-1686564889672).
SparkSession available as 'spark'.
>>> █
```

## spark-shell

## historyserver



## Przykładowy skrypt (pyspark + pyspark.sql)

```python
from pyspark.sql import SparkSession

covid = "/datasets/covid-dataset.jsonl"
spark = SparkSession.builder.appName("CovidApp").getOrCreate()

df = spark.read.json(covid)
df.select("location").write.csv('/spark-result/dataframe-select',
header=True)

spark.read.json(covid).createOrReplaceTempView("covid")
spark.sql("SELECT location FROM covid").write.csv('/spark-result/sql-
select', header=True)

spark.stop()
```

## Uruchomienie z poziomu jupyter notebook

```
_ = run_in_master("spark-submit --master yarn --deploy-mode cluster
/data/master_volume/spark_scripts/test.py")
_ = merge_results("/spark-result/sql-select")
print_hdfs_output("/spark-result/sql-select")
```

# Uruchomiona aplikacja na hadoopie