

**EKONOMICKÁ UNIVERZITA V BRATISLAVE  
FAKULTA HOSPODÁRSKEJ INFORMATIKY**

**SPRACOVANIE A ANALÝZA LOGOV V  
NADSTAVBE HIVE V SYSTÉME LINUX**

**Zadanie z predmetu Big Data**

**2023**

**Bc. Ondrej Šima  
Bc. Alena Stracenská**

# Obsah

|   |           |
|---|-----------|
| <b>Úvod</b>   | <b>4</b>  |
| <b>1 Nadstavba Hive na Linuxe</b>                   | <b>5</b>  |
| 1.1 Inštalácia . . . . .                            | 5         |
| 1.2 Inštalácia databázy Apache Derby . . . . .      | 6         |
| 1.3 Konfigurácia Apache Hive . . . . .              | 6         |
| <b>2 Práca s logmi v Hive</b>                       | <b>9</b>  |
| 2.1 Spracovanie a analýza logov . . . . .           | 10        |
| 2.1.1 Spracovanie a analýza 3.2 GB súboru . . . . . | 13        |
| <b>Záver</b>  | <b>15</b> |
| <b>Zoznam použitej literatúry</b>                   | <b>16</b> |

## Zoznam obrázkov a tabuliek

|           |  |    |
|-----------|--|----|
| Obrázok 1 | Ukážka súboru s logmi, zdroj: [vlastné spracovanie] . . . . .                              | 9  |
| Obrázok 2 | Ukážka vytvorenia tabuľky a načítania dát, zdroj: [vlastné spracovanie] . . . . .          | 11 |
| Obrázok 3 | Proces vykonávania jobs pri agregácii, zdroj: [vlastné spracovanie]                        | 12 |
| Obrázok 4 | Výsledok SQL dopytu agregácie IP adries, zdroj: [vlastné spracovanie] . . . . .            | 12 |
| Obrázok 5 | Vytvorenie tabuľky a načítanie dát z celého súboru, zdroj: [vlastné spracovanie] . . . . . | 13 |
| Obrázok 6 | Proces vykonávania jobs pri agregácii, zdroj: [vlastné spracovanie]                        | 13 |
| Obrázok 7 | Výsledok SQL dopytu agregácie IP adries, zdroj: [vlastné spracovanie] . . . . .            | 14 |

# Úvod

V tomto zadaní si ukážeme ako správne nainštalovať nadstavbu Hive na operačný systém Linux. Celá inštalácia pozostáva z niekoľkých krokov, ktoré sú síce zdĺhavejšie, ale vo výsledku nám prinesú funkčnú nadstavbu Hadoopu.

V ďalšom kroku si predstavíme použité logy, ich zdroj a štruktúru. Následne si ukážeme ako je možné logy spracovať cez Hive, od niekoľko megabajtového súboru až po niekoľko gigabajtový súbor. Potom si ukážeme agregácie, pomocou ktorých vieme zistiť z logov zaujímavé informácie.

Veríme, že tento dokument bude prínosný pre čitateľa nie len po teoretickej stránke, ale aj po praktickej, kde sa naučí čo všetko inštalácia Hive na Linux obnáša a ako je možné spracovať a analyzovať pomocou tejto populárnej nadstavby logy.

# 1 Nadstavba Hive na Linuxe

Keďže z predchádzajúcich zadaní sme disponovali funkčnou inštaláciou Hadoopu, kroky v podobe overenia verzie Javy, jej inštaláciu, konfiguráciu a overenie funkčnosti samotného Hadoopu sme preskočili. Pri práci s Apache Hive pokračujeme na rovnakom systéme z minulého zadania. Pred prácou s Hive si zapneme potrebné Hadoop služby (HDFS a YARN).

## 1.1 Inštalácia

Verziu Hive sme vybrali podľa nainštalovanej verzie Hadoopu, v našom prípade pre Hadoop 3.3.0 stiahneme Hive verziu 3.4.3. Stiahnutý balík následne rozbalíme pomocou:

```
$ tar zxvf apache-hive-0.14.0-bin.tar.gz
```

Úspešné rozbalenie môžeme overiť existenciou rozbaleného priečinka so súbormi Hive:

```
$ ls | grep hive
apache-hive-3.1.3-bin.tar.gz
archive.zip
```

Priečinkok obsahujúci binárne a konfiguračné súbory Hive skopírujeme do /usr/local/hive:

```
# mv apache-hive-3.1.3-bin /usr/local/hive
```

Podobne ako pri inštalácii Hadoopu si vytvoríme systémové premenné pre cesty ku Hive súborom, tie vložíme na koniec súboru bashrc:

```
$ echo "export HIVE_HOME=/usr/local/hive" >> ~/.bashrc
$ echo "export PATH=$PATH:$HIVE_HOME/bin" >> ~/.bashrc
$ echo "export CLASSPATH=$CLASSPATH:/usr/local/Hadoop/lib/*:." >> ~/.bashrc
$ echo "export CLASSPATH=$CLASSPATH:/usr/local/hive/lib/*:." >> ~/.bashrc
$ source ~/.bashrc
```

Do priečinka /usr/local/hive sa teraz môžeme zjednodušene dostať pomocou premennej \$HIVE\_HOME a pokračovať konfiguráciou.

## 1.2 Inštalácia databázy Apache Derby

Pre prácu s Hive potrebujeme tiež nainštalovať databázový server, pomocou ktorého si Hive vytvorí spravuje a číta Metastore databázu. Použiť je možné akúkoľvek SQL databázu, napríklad MySQL alebo MariaDB. V našom prípade sme sa rozhodli použiť databázu Apache Derby.

Po stiahnutí Derby databázy verzie 10.4.2.0 sme rozbalili archív pomocou:

```
$ tar zxvf db-derby-10.4.2.0-bin.tar.gz
```

Priečinok s binárnymi súbormi sme skopírovali do /usr/local:

```
# mv db-derby-10.4.2.0-bin /usr/local/derby
```

Podobne ako v prípade Hive, nastavili sme si potrebné systémové premenné:

```
$ echo "export DERBY_HOME=/usr/local/derby" >> ~/.bashrc
$ echo "export PATH=$PATH:$DERBY_HOME/bin" >> ~/.bashrc
$ echo "export CLASSPATH=$CLASSPATH:$DERBY_HOME/lib/derby.jar:$DERBY_HOME/lib/derbytools.jar" >> ~/.bashrc
$ source ~/.bashrc
```

## 1.3 Konfigurácia Apache Hive

Hive používa podobnú schému podpriechinkov ako Hadoop v predošlom zadaní, konfiguračné súbory nájdeme v priečinku \$HIVE\_HOME/conf. Tu nájdeme predpripravený vzor konfigurácie a shellsript, ktoré sú neaktívne (majú nevhodnú príponu), môžeme ich premenovať alebo skopírovať so správnym názvom:

```
$ cd $HIVE_HOME/conf
$ cp hive-env.sh.template hive-env.sh
$ cp hive-default.xml.template hive-site.xml
```

V súbore hive-env.sh doplníme cestu ku Hadoopu:

```
export HADOOP_HOME=/usr/local/hadoop
```

V konfiguračnom súbore hive-site.xml sme potom vykonali niekoľko zmien. Tieto zmeny sme vykonávali postupne, tak ako sme narážali na dané problémy, kvôli zjednodušeniu procesu opíšeme všetky zmeny naraz.

Vo verzii Hive, ktorý sme stiahli je chyba v konfiguračnom súbore hive-env.sh. Na riadku 3220 je použitý neplatný znak, z dôvodu čoho pravdepodobne dochádza ku zlému čítaniu tohto súboru. Chyba sa nachádza na riadku 3220, textové editory ako napríklad Nano ale Vim tento znak zvýraznia. Ďalej sme do tohto konfiguračného súboru doplnili nasledujúce:

```
<property>
  <name>system:java.io.tmpdir</name>
  <value>/tmp/hive/java</value>
</property>
<property>
  <name>system:user.name</name>
  <value>${user.name}</value>
</property>
<property>
  <name>javax.jdo.option.ConnectionURL</name>
  <value>jdbc:derby://localhost:1527/metastore_db;create=true </value>
  <description>JDBC connect string for a JDBC metastore </description>
</property>
<property>
  <name>javax.jdo.option.ConnectionDriverName</name>
  <value>org.apache.derby.jdbc.ClientDriver</value>
  <description>Driver class name for a JDBC metastore</description>
</property>
```

Potom sme v priečinku conf vytvorili súbor s názvom jpox.properties a vložili do neho nasledujúci obsah:

```
javax.jdo.PersistenceManagerFactoryClass =

org.jpox.PersistenceManagerFactoryImpl
org.jpox.autoCreateSchema = false
org.jpox.validateTables = false
org.jpox.validateColumns = false
org.jpox.validateConstraints = false
org.jpox.storeManagerType = rdbms
org.jpox.autoCreateSchema = true
org.jpox.autoStartMechanismMode = checked
org.jpox.transactionIsolation = read_committed
```

```
javax.jdo.option.DetachAllOnCommit = true
javax.jdo.option.NontransactionalRead = true
javax.jdo.option.ConnectionDriverName = org.apache.derby.jdbc.ClientDriver
javax.jdo.option.ConnectionURL = jdbc:derby://hadoop1:1527/metastore_db;create = true
javax.jdo.option.ConnectionUserName = APP
javax.jdo.option.ConnectionPassword = mine
```

Jeden z ďalších problémov, na ktorý sme narazili je inkonzistencia verzií knižnice guava.jar medzi inštaláciou Hive a Hadoop. Zatiaľ čo Hive používal verziu 19.0, Hadoop mal k dispozícii novšiu verziu, 27.0, z dôvodu čoho Hive nefungoval. Tento problém sme vyriešili nahradením staršej verzie v priečinku `$HIVE_HOME/lib` novšou verziou z Hive:

```
$ sudo mv /usr/local/hadoop/share/hadoop/hdfs/lib/guava-27.0-jre.jar
/usr/local/hive/lib/guava-27.0-jre.jar
```

Pred samotnou prácou v Hive si vytvoríme Metastore databázu v ktorej si Hive ukladá informácie, ako schémy, ku samotným dátam v Hadoope. Môžeme tak urobiť pomocou:

```
$ cd $HIVE_HOME
$ hive --service metastore
```

Nakoniec si pre Hive vytvoríme v HDFS priečinky a priradíme im práva:

```
$ $HADOOP_HOME/bin/hadoop fs -mkdir /tmp
$ $HADOOP_HOME/bin/hadoop fs -mkdir -p /user/hive/warehouse
$ $HADOOP_HOME/bin/hadoop fs -chmod g+w /tmp
$ $HADOOP_HOME/bin/hadoop fs -chmod g+w /user/hive/warehouse
```

Funkčnosť inštalácie môžeme overiť správnym spustením Hive a úspešným zadaním príkazu, napríklad výpisom všetkých tabuliek:

```
$ cd $HIVE_HOME
$ hive
>show tables;
```

Ak je Hive nainštalovaný správne mal by zobrazit:

```
OK
Time taken: 2.798 seconds
```



Alternatívne, ak plánujeme Hive využívať ako server, môžeme ho spustiť pomocou príkazu:

```
$ hive --service hiveserver2 --hiveconf hive.server2.thrift.port=10000
--hiveconf hive.root.logger=INFO,console
```

Potom sa na Hive môžeme pripojiť napríklad pomocou JDBC konektora z Java alebo Python aplikácie (aktuálna verzia balíčka implementujúceho JDBC konektor pre Python nejde nainštalovať z dôvodu závislosti na už nedostupnej systémovej knižnici libsasl2-dev). Alternatívne môžeme funkčnosť Hive servera overiť pomocou nástroja beeline:

```
$ beeline -u jdbc:hive2://localhost:10000/
```

V prípade funkčného pripojenia sa nám zobrazí hive prompt a bude možné vykonávať príkazy. [1]

## 2 Práca s logmi v Hive

Dáta, s ktorými sme pracovali boli vo forme apačovských webových logov. Môžeme si ich stiahnuť na kaggle.com. Pre naše potreby sme po stiahnutí dát vytiahli 100 000 riadkov, celý súbor bude spracovaný v kapitole 2.1.1. Sú to pološtruktúrované dáta znázornené na obrázku 1.



Obrázok 1: Ukážka súboru s logmi, zdroj: [vlastné spracovanie]

## 2.1 Spracovanie a analýza logov

Základom pre spracovanie a analýzu logov bolo vytvoriť si tabuľku, do ktorej sme pomocou regexu parsovali jednotlivé logy. Tabuľku sme rozdelili na 10 polí a pomenovali podľa jednotlivých zložiek. Ak sme tabuľku vopred už mali vytvorenú, tak sme si ju vymazali.

```
DROP TABLE IF EXISTS access_100K PURGE;
```

```
CREATE TABLE access_100k (  
  ip STRING,  
  datetime STRING,  
  method STRING,  
  uri STRING,  
  http_version STRING,  
  http_response STRING,  
  response_size STRING,  
  full_url STRING,  
  user_agent STRING,  
  other STRING  
)
```

```
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.RegexSerDe'  
WITH SERDEPROPERTIES (  
  "input.regex" = "^(\\S*).*\\[(.*)\\]\\s\"(\\S*)\\s(\\S*)  
  \\s([^\"]*)\"\\s(\\S*)\\s(\\S*)\\s\"([^\"]*)\"\\s\"([^\"]*)\"\\s\"([^\"]*)\""  
)
```

Regex pre členenie logu sme si spravili na <https://regex101.com/r/z9JTRm/1>. Po otvorení si vieme pozrieť pekne jednotlivé vyparsované položky logu a pochopiť syntax regexu. Následne sme vyššie uvedený kód vložili do shellu. Výsledkom je vytvorená tabuľka, do ktorej sme následne cez load data načítali logy pomocou príkazu:

```
LOAD DATA LOCAL INPATH "/home/alena/access_100K.log" INTO TABLE access_100k;
```



```
Activities Terminal mar 29 20:44 alena@alena-virtual-machine: /usr/local/hive

In order to limit the maximum number of reducers:
set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Starting Job = job_1608109839198_0009, Tracking URL = http://alena-virtual-machine:8088/proxy/application_1608109839198_0009/
Kill Command = /usr/local/hadoop/bin/mapred job -kill job_1608109839198_0009
Hadoop Job Information for Stage-1: number of mappers: 1; number of reducers: 1
2023-03-29 20:42:03,922 Stage-1 map = 0%, reduce = 0%
2023-03-29 20:42:16,426 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 9.26 sec
2023-03-29 20:42:22,399 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 11.17 sec
MapReduce Total cumulative CPU time: 11 seconds 370 msec
Ended Job = job_1608109839198_0009
Launching Job 2 out of 3
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Starting Job = job_1608109839198_0010, Tracking URL = http://alena-virtual-machine:8088/proxy/application_1608109839198_0010/
Kill Command = /usr/local/hadoop/bin/mapred job -kill job_1608109839198_0010
Hadoop Job Information for Stage-2: number of mappers: 1; number of reducers: 1
2023-03-29 20:42:36,690 Stage-2 map = 0%, reduce = 0%
2023-03-29 20:42:42,061 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 1.65 sec
2023-03-29 20:42:48,199 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 3.2 sec
MapReduce Total cumulative CPU time: 3 seconds 200 msec
Ended Job = job_1608109839198_0010
Launching Job 3 out of 3
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Starting Job = job_1608109839198_0011, Tracking URL = http://alena-virtual-machine:8088/proxy/application_1608109839198_0011/
Kill Command = /usr/local/hadoop/bin/mapred job -kill job_1608109839198_0011
Hadoop Job Information for Stage-3: number of mappers: 1; number of reducers: 1
2023-03-29 20:43:00,657 Stage-3 map = 0%, reduce = 0%
2023-03-29 20:43:04,785 Stage-3 map = 100%, reduce = 0%, Cumulative CPU 1.4 sec
2023-03-29 20:43:09,961 Stage-3 map = 100%, reduce = 100%, Cumulative CPU 3.29 sec
MapReduce Total cumulative CPU time: 3 seconds 290 msec
Ended Job = job_1608109839198_0011
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 11.17 sec HDFS Read: 31130838 HDFS Write: 126275 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 3.2 sec HDFS Read: 132971 HDFS Write: 394 SUCCESS
Stage-Stage-3: Map: 1 Reduce: 1 Cumulative CPU: 3.29 sec HDFS Read: 7829 HDFS Write: 357 SUCCESS
Total MapReduce CPU Time Spent: 17 seconds 660 msec
```

Obrázok 3: Proces vykonávania jobs pri agregácii, zdroj: [vlastné spracovanie]

```
Activities Terminal mar 29 20:44 alena@alena-virtual-machine: /usr/local/hive

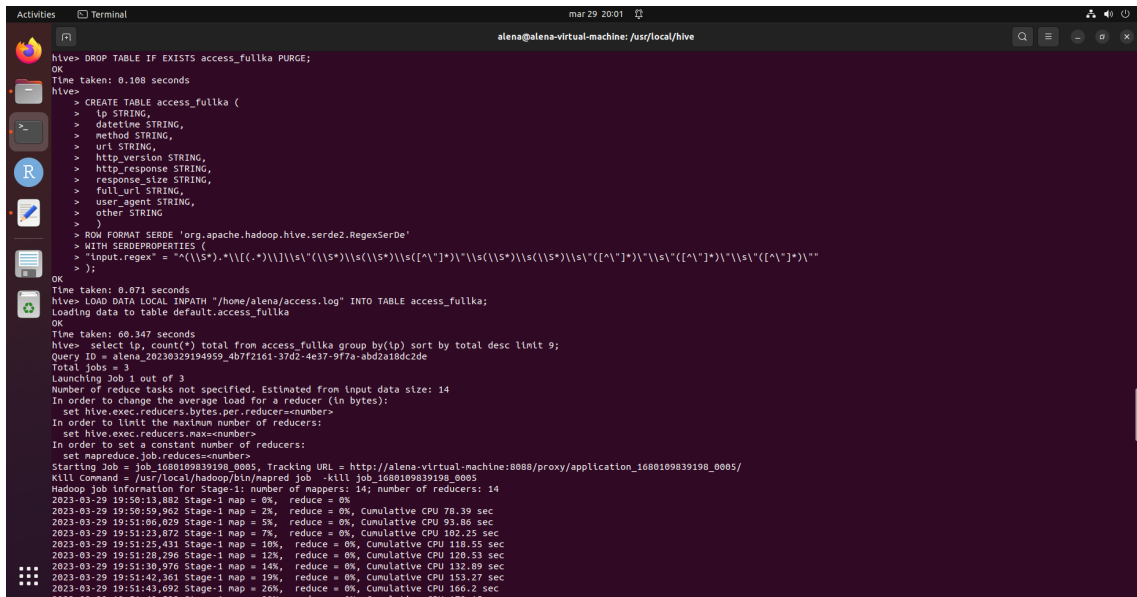
Launching Job 2 out of 3
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Starting Job = job_1608109839198_0010, Tracking URL = http://alena-virtual-machine:8088/proxy/application_1608109839198_0010/
Kill Command = /usr/local/hadoop/bin/mapred job -kill job_1608109839198_0010
Hadoop Job Information for Stage-2: number of mappers: 1; number of reducers: 1
2023-03-29 20:42:36,690 Stage-2 map = 0%, reduce = 0%
2023-03-29 20:42:42,061 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 1.65 sec
2023-03-29 20:42:48,199 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 3.2 sec
MapReduce Total cumulative CPU time: 3 seconds 200 msec
Ended Job = job_1608109839198_0010
Launching Job 3 out of 3
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Starting Job = job_1608109839198_0011, Tracking URL = http://alena-virtual-machine:8088/proxy/application_1608109839198_0011/
Kill Command = /usr/local/hadoop/bin/mapred job -kill job_1608109839198_0011
Hadoop Job Information for Stage-3: number of mappers: 1; number of reducers: 1
2023-03-29 20:43:00,657 Stage-3 map = 0%, reduce = 0%
2023-03-29 20:43:04,785 Stage-3 map = 100%, reduce = 0%, Cumulative CPU 1.4 sec
2023-03-29 20:43:09,961 Stage-3 map = 100%, reduce = 100%, Cumulative CPU 3.29 sec
MapReduce Total cumulative CPU time: 3 seconds 290 msec
Ended Job = job_1608109839198_0011
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 11.17 sec HDFS Read: 31130838 HDFS Write: 126275 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 3.2 sec HDFS Read: 132971 HDFS Write: 394 SUCCESS
Stage-Stage-3: Map: 1 Reduce: 1 Cumulative CPU: 3.29 sec HDFS Read: 7829 HDFS Write: 357 SUCCESS
Total MapReduce CPU Time Spent: 17 seconds 660 msec
OK
66.249.66.194 9331
66.249.66.91 6202
23.101.169.3 2863
91.90.72.15 2398
40.77.167.170 2292
207.46.13.9 1915
5.237.16.117 1132
207.46.13.136 882
Time taken: 79.117 seconds, Fetched: 9 row(s)
hive:
```

Obrázok 4: Výsledok SQL dopytu agregácie IP adries, zdroj: [vlastné spracovanie]

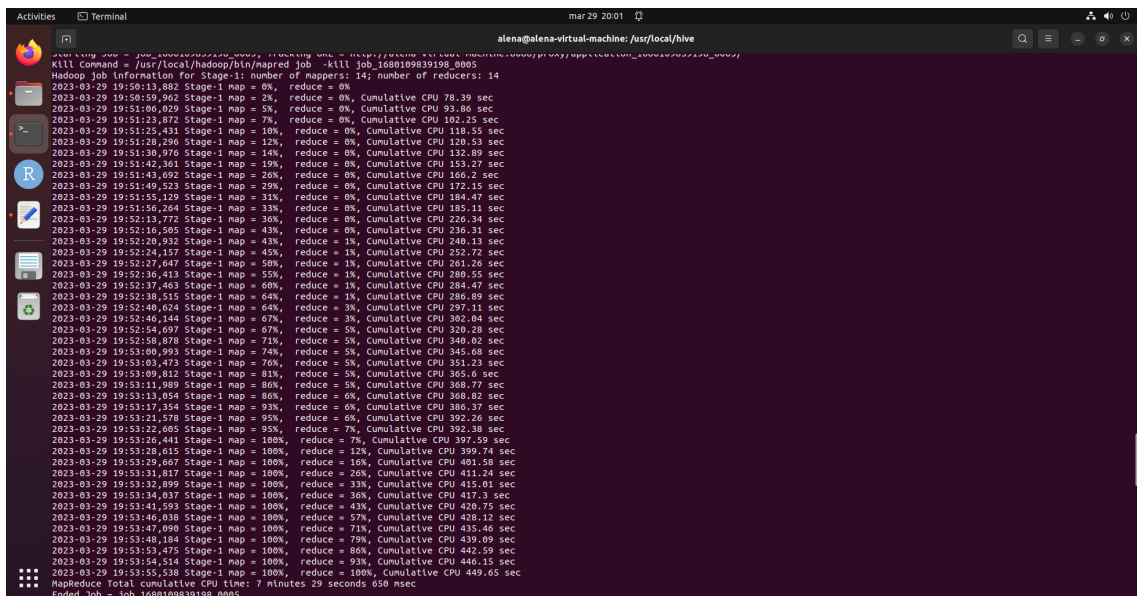
Podobným spôsobom by sme pomocou ďalšieho SQL dopytu mohli napríklad rozdeliť čas na jednotlivé zložky, ak by sme chceli zobraziť grafy prípadne iné metriky.

### 2.1.1 Spracovanie a analýza 3.2 GB súboru

To isté sme spravili taktiež pre celý súbor s logmi, ktorý má veľkosť 3.2 GB. V Hortonworkse sme sa pri načítavaní takto veľkého súboru stretli s dlhým časovým trvaním, z tohto dôvodu je ukážka spracovania celého súboru len v tomto zadaní.



**Obrázok 5:** Vytvorenie tabuľky a načítanie dát z celého súboru, zdroj: [vlastné spracovanie]



**Obrázok 6:** Proces vykonávania jobs pri agregácii, zdroj: [vlastné spracovanie]

```
Activities Terminal mar 29 20:01 alena@alena-virtual-machine: /usr/local/hive

Launching Job 2 out of 3
Number of reduce tasks not specified. Estimated from input data size: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1680109839198_0006, Tracking URL = http://alena-virtual-machine:8088/proxy/application_1680109839198_0006/
Kill Command = /usr/local/hadoop/bin/mapred job -kill job_1680109839198_0006
Hadoop job information for stage-2: number of mappers: 1; number of reducers: 1
2023-03-29 19:54:09,674 Stage-2 map = 0%, reduce = 0%
2023-03-29 19:54:18,102 Stage-2 map = 100%, reduce = 0%, Cumulative CPU 4.66 sec
2023-03-29 19:54:22,239 Stage-2 map = 100%, reduce = 100%, Cumulative CPU 6.45 sec
MapReduce Total cumulative CPU time: 6 seconds 450 msec
Ended Job = job_1680109839198_0006
Launching Job 3 out of 3
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Starting Job = job_1680109839198_0007, Tracking URL = http://alena-virtual-machine:8088/proxy/application_1680109839198_0007/
Kill Command = /usr/local/hadoop/bin/mapred job -kill job_1680109839198_0007
Hadoop job information for stage-3: number of mappers: 1; number of reducers: 1
2023-03-29 19:54:37,888 Stage-3 map = 0%, reduce = 0%
2023-03-29 19:54:42,049 Stage-3 map = 100%, reduce = 0%, Cumulative CPU 1.44 sec
2023-03-29 19:54:47,108 Stage-3 map = 100%, reduce = 100%, Cumulative CPU 3.36 sec
MapReduce Total cumulative CPU time: 3 seconds 360 msec
Ended Job = job_1680109839198_0007
MapReduce Jobs Launched:
Stage-Stage-1: Map: 14 Reduce: 14 Cumulative CPU: 449.65 sec HDFS Read: 3502736782 HDFS Write: 8193626 SUCCESS
Stage-Stage-2: Map: 1 Reduce: 1 Cumulative CPU: 6.45 sec HDFS Read: 8203663 HDFS Write: 399 SUCCESS
Stage-Stage-3: Map: 1 Reduce: 1 Cumulative CPU: 3.36 sec HDFS Read: 7834 HDFS Write: 370 SUCCESS
Total MapReduce CPU Time Spent: 7 minutes 39 seconds 460 msec
OK
66.249.66.194 353483
66.249.66.91 314522
151.239.241.163 92473
66.249.66.92 88332
91.99.38.32 45973
104.222.32.91 42858
91.99.72.15 38694
91.99.47.57 38689
5.78.190.233 37203
Time taken: 289.21 seconds, Fetched: 9 row(s)
hive> select count(*) from access_fulika;
```

Obrázok 7: Výsledok SQL dopytu agregácie IP adries, zdroj: [vlastné spracovanie]

## Záver

Cieľom tohto zadania bolo opísať inštaláciu Hive na Linuxe a následne spracovať a analyzovať niekoľkomegabajtový súbor až po trojgigabajtový súbor. Počas inštalácie sme sa stretli s problémami, ktoré sme opísali vyššie v dokumente, ale aj napriek tomu sme dokázali náš cieľ spracovania a analýzy logov splniť.

V porovnaní s produktom Hortonworks Sandbox je spracovanie logov v systéme Linux intuitívnejšie a rýchlejšie. Taktiež je veľmi natívne a povedali by sme aj jednoduchšie oproti Hortonworks Sandbox, keďže už spomenutý Hortonworks nám prišiel pomalší a zložitejší na naučenie, taktiež sme v ňom mali problém s administrátorskými právami a načítaním niekoľkogigabajtového súboru. Preto by sme na základe porovnania Hive v Linuxe a v Hortonworkse odporučili použiť Hive natívne pod Linuxom.

Veríme, že toto zadanie prinesie čitateľovi nový pohľad na možnosť spracovávať veľké súbory a získa nové prípadne obohatí už doterajšie vedomosti a schopnosti o danej problematike, ktorá je tiež súčasťou Big Data.

## Zoznam použitej literatúry

- [1] TUTORIALS POINT INDIA PRIVATE LIMITED. 2023. *Hive - Installation* In tutorialspoint.com [online]. 2023. [citované dňa 29.03.2023]. Dostupné na internete: [https://www.tutorialspoint.com/hive/hive\\_installation.htm](https://www.tutorialspoint.com/hive/hive_installation.htm).