

Создание кластера из виртуальных машин

Проект по информатике

Стафеев И.А., Голованов Д.И., Шарыпов Е.А.

Санкт-Петербург, 2023

О чем лабораторная

Цель: развертывание кластера из виртуальных машин с установленной на них распределенной файловой системой, чтобы созданный кластер мог служить средством обработки больших данных в соответствии с парадигмой распределенных вычислений MapReduce.

Задачи

1. Придумать практическую задачу, решение которой может быть осуществлено с помощью кластера
2. Создать и настроить ВМ для установки РФС
3. Установить РФС
4. Написать скрипты для MapReduce-задачи
5. Провести несколько раз обработку входных данных, доказав работоспособность кластера

Почему именно такая тема и цель

- Охватывает множество тем из изученного курса информатики
- Это интересно и нескучно
- У нас имеется небольшой опыт работы в этой сфере

Немного теории

MapReduce - модель распределённых вычислений. Map – предварительная обработка данных, Reduce – свертка данных

Apache Hadoop - это инструмент для распределенного хранения данных (HDFS + YARN)

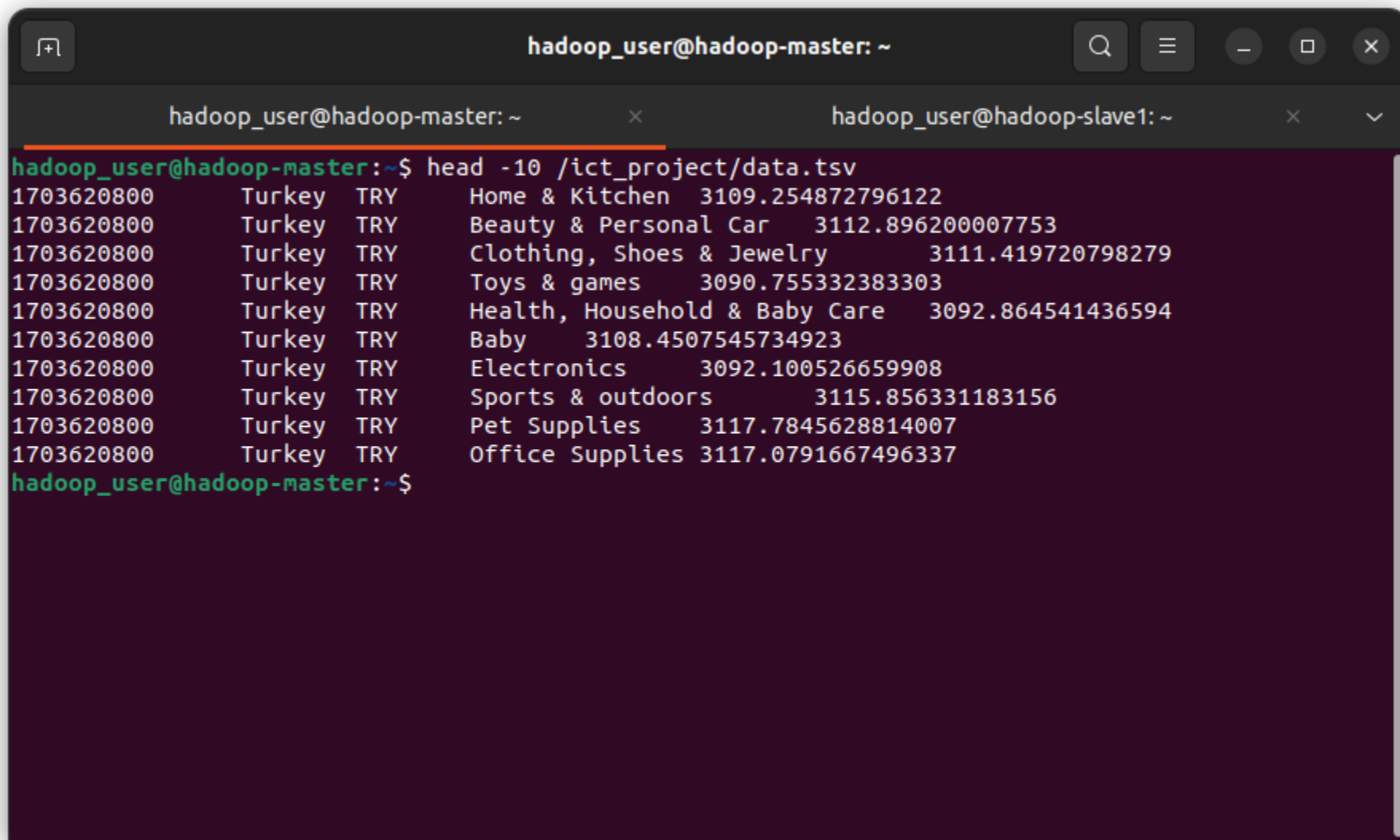
Apache YARN - система планирования заданий
и управления кластером

Apache Airflow - инструмент для оркестрации данных и построения **DAG**'ов – представления задач работы в виде графа

Практическая задача

Преобразование данных о покупках. Раз в 20 минут подсчитывать выручку в рублях и строить график распределения выручки по странам для категории товаров с наибольшей выручкой

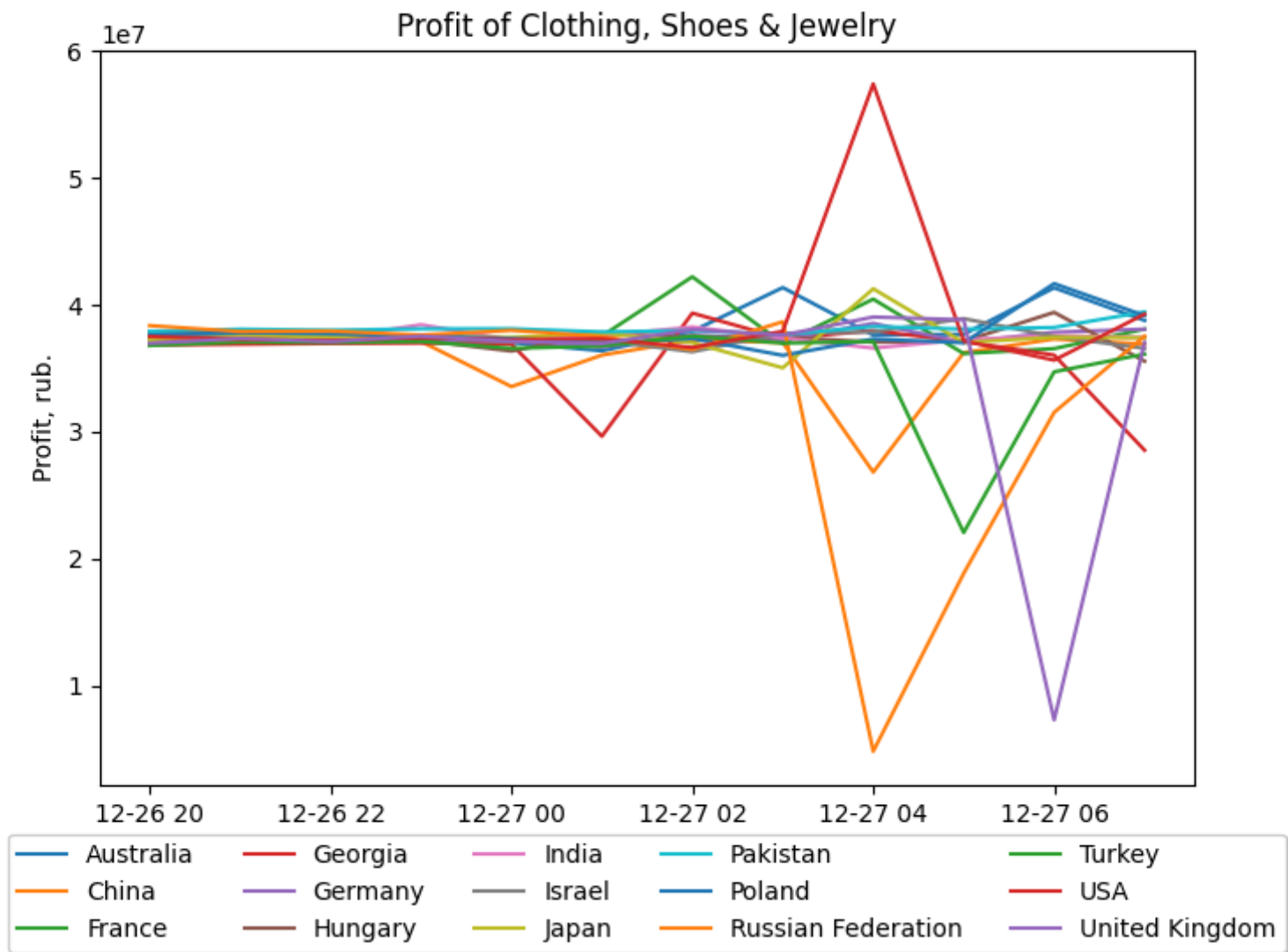
Dataset



A terminal window titled 'hadoop_user@hadoop-master: ~' with standard window controls. It shows a command prompt where the user has run 'head -10 /ict_project/data.tsv'. The output displays ten rows of data, each with a numeric ID, a country, a category, and a float value. The window also shows a second tab for 'hadoop_user@hadoop-slave1: ~'.

```
hadoop_user@hadoop-master: ~  
hadoop_user@hadoop-master:~$ head -10 /ict_project/data.tsv  
1703620800      Turkey TRY      Home & Kitchen  3109.254872796122  
1703620800      Turkey TRY      Beauty & Personal Car  3112.896200007753  
1703620800      Turkey TRY      Clothing, Shoes & Jewelry  3111.419720798279  
1703620800      Turkey TRY      Toys & games  3090.755332383303  
1703620800      Turkey TRY      Health, Household & Baby Care  3092.864541436594  
1703620800      Turkey TRY      Baby  3108.4507545734923  
1703620800      Turkey TRY      Electronics  3092.100526659908  
1703620800      Turkey TRY      Sports & outdoors  3115.856331183156  
1703620800      Turkey TRY      Pet Supplies  3117.7845628814007  
1703620800      Turkey TRY      Office Supplies  3117.0791667496337  
hadoop_user@hadoop-master:~$
```

ХОТИМ



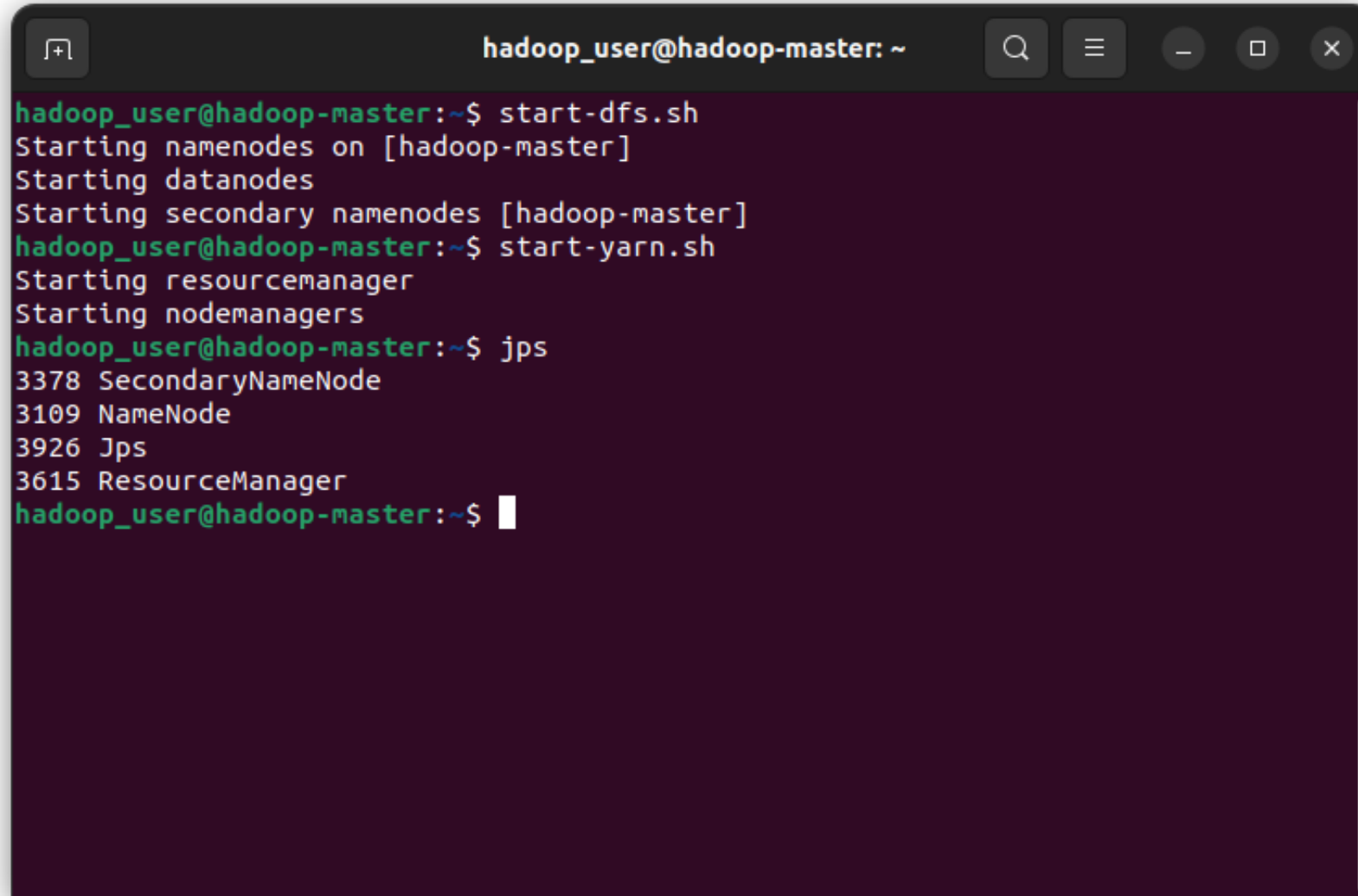
Этап 0

- Создание трех VM и подключение в одну сеть NAT
- Создание между ними SSH-подключения
- Установка Java

Этап 1

- Установка и конфигурация HDFS
- Форматирование и запуск распределенной файловой системы на кластере

Запуск HDFS



```
hadoop_user@hadoop-master: ~  
hadoop_user@hadoop-master:~$ start-dfs.sh  
Starting namenodes on [hadoop-master]  
Starting datanodes  
Starting secondary namenodes [hadoop-master]  
hadoop_user@hadoop-master:~$ start-yarn.sh  
Starting resourcemanager  
Starting nodemanagers  
hadoop_user@hadoop-master:~$ jps  
3378 SecondaryNameNode  
3109 NameNode  
3926 Jps  
3615 ResourceManager  
hadoop_user@hadoop-master:~$
```

Рабочие ноды кластера

Namenode information

hdfs:///user/student/inp...

localhost:9870/dfshealth.html#tab-datanode

In operation

DataNode State: All

Show: 25 entries

Search:

Node	Http Address	Last contact	Last Block Report	Used	Non DFS Used	Capacity	Blocks	Block pool used	Version	
✓/default-rack/hadoop-slave1:9866 (10.0.2.9:9866)	http://hadoop-slave1:9864	1s	1m	24 KB	6.63 GB	9.98 GB	<div><div></div></div>	0	24 KB (0%)	3.3.6
✓/default-rack/hadoop-slave2:9866 (10.0.2.10:9866)	http://hadoop-slave2:9864	2s	1m	24 KB	6.61 GB	9.98 GB	<div><div></div></div>	0	24 KB (0%)	3.3.6

Showing 1 to 2 of 2 entries

Previous

1

Next

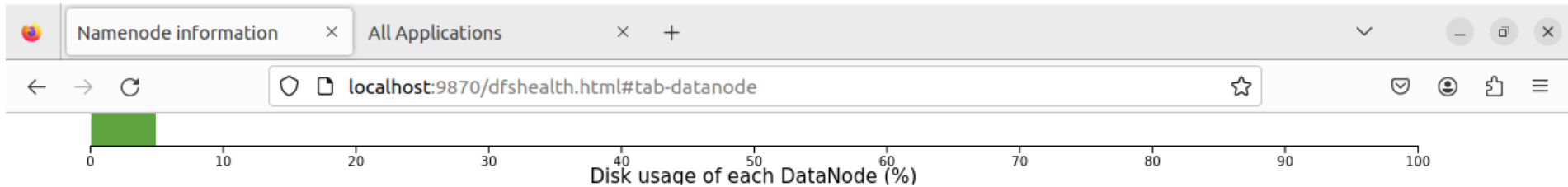
Entering Maintenance

No nodes are entering maintenance.

Decommissioning

Этап 2

- Создание скрипта для получения валют
- Создание скрипта с map и reduce функциями
- Создание скрипта для визуализации
- Последовательный запуск скриптов над кластерах



In operation

DataNode State: All Show 25 entries Search:

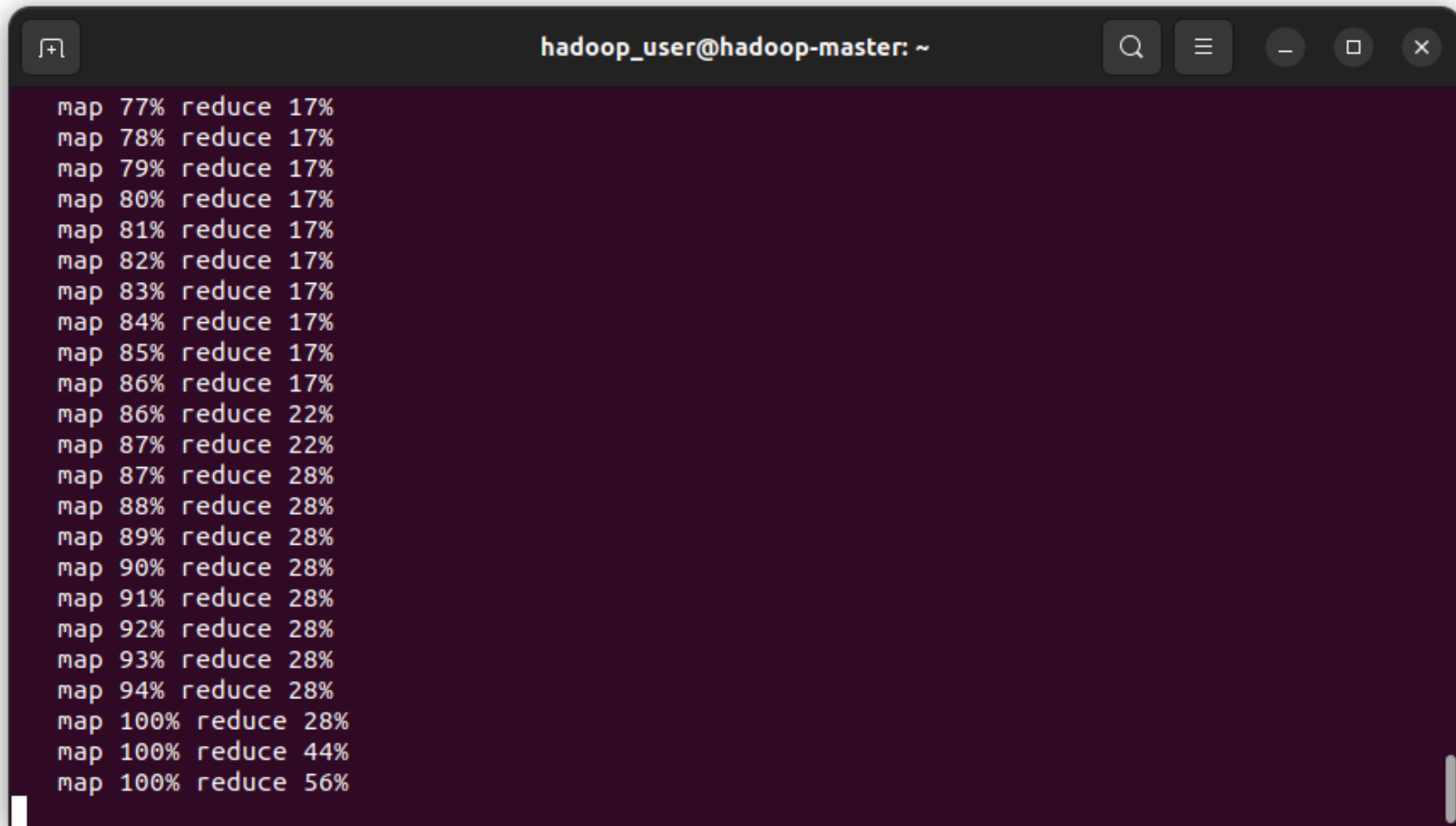
Node	Http Address	Last contact	Last Block Report	Used	Non DFS Used	Capacity	Blocks	Block pool used	Version
✓ /default-rack/hadoop-slave1:9866 (10.0.2.9:9866)	http://hadoop-slave1:9864	0s	27m	261.66 MB	6.77 GB	9.98 GB	127	261.66 MB (2.56%)	3.3.6
✓ /default-rack/hadoop-slave2:9866 (10.0.2.10:9866)	http://hadoop-slave2:9864	0s	27m	148.66 MB	6.78 GB	9.98 GB	138	148.66 MB (1.45%)	3.3.6

Showing 1 to 2 of 2 entries

Previous 1 Next

Entering Maintenance

```
hadoop_user@hadoop-master: ~  
hadoop_user@hadoop-master:~$ lab/get_currencies.sh  
% Total      % Received % Xferd  Average Speed   Time    Time     Time  Current  
           Dload  Upload   Total   Spent    Left   Speed  
100 2844    0 2844    0     0 10472      0 --:--:-- --:--:-- --:--:-- 10455  
hadoop_user@hadoop-master:~$ python3 lab/mapred_script.py --currencies lab/currencies.json -r  
hadoop hdfs:///user/student/input/data.tsv > lab/output.tsv  
No configs found; falling back on auto-configuration  
No configs specified for hadoop runner  
Looking for hadoop binary in /usr/local/hadoop/bin...  
Found hadoop binary: /usr/local/hadoop/bin/hadoop  
Using Hadoop version 3.3.6  
Looking for Hadoop streaming jar in /usr/local/hadoop...  
Found Hadoop streaming jar: /usr/local/hadoop/share/hadoop/tools/lib/hadoop-streaming-3.3.6.jar  
Creating temp directory /tmp/mapred_script.hadoop_user.20231226.204537.922149  
uploading working dir files to hdfs:///user/hadoop_user/tmp/mrjob/mapred_script.hadoop_user.20231226.204537.922149/files/wd...  
Copying other local files to hdfs:///user/hadoop_user/tmp/mrjob/mapred_script.hadoop_user.20231226.204537.922149/files/  
Running step 1 of 1...  
packageJobJar: [/tmp/hadoop-unjar2683409111838308545/] [] /tmp/streamjob628156531139366329.jar tmpDir=null
```

A terminal window with a dark background and light text. The title bar at the top shows the user 'hadoop_user' on a machine named 'hadoop-master' in the home directory. On the right side of the title bar are standard window controls: a search icon, a menu icon, and buttons for minimize, maximize, and close. The terminal content displays the progress of a Hadoop job, with 'map' and 'reduce' tasks and their respective completion percentages. The progress increases from 77% map / 17% reduce to 100% map / 56% reduce. A vertical scrollbar is visible on the right edge of the terminal area.

```
hadoop_user@hadoop-master: ~  
map 77% reduce 17%  
map 78% reduce 17%  
map 79% reduce 17%  
map 80% reduce 17%  
map 81% reduce 17%  
map 82% reduce 17%  
map 83% reduce 17%  
map 84% reduce 17%  
map 85% reduce 17%  
map 86% reduce 17%  
map 86% reduce 22%  
map 87% reduce 22%  
map 87% reduce 28%  
map 88% reduce 28%  
map 89% reduce 28%  
map 90% reduce 28%  
map 91% reduce 28%  
map 92% reduce 28%  
map 93% reduce 28%  
map 94% reduce 28%  
map 100% reduce 28%  
map 100% reduce 44%  
map 100% reduce 56%
```



```
hadoop_user@hadoop-master: ~  
CONNECTION=0  
IO_ERROR=0  
WRONG_LENGTH=0  
WRONG_MAP=0  
WRONG_REDUCE=0  
job output is in hdfs:///user/hadoop_user/tmp/mrjob/mapred_script.hadoop_user.20231226.204537  
.922149/output  
Streaming final output from hdfs:///user/hadoop_user/tmp/mrjob/mapred_script.hadoop_user.2023  
1226.204537.922149/output...  
Removing HDFS temp directory hdfs:///user/hadoop_user/tmp/mrjob/mapred_script.hadoop_user.202  
31226.204537.922149...  
Removing temp directory /tmp/mapred_script.hadoop_user.20231226.204537.922149...  
hadoop_user@hadoop-master:~$ head -10 lab/output.tsv  
["2023-12-26T20", "Australia", "Baby"] 37507092.75928125  
["2023-12-26T20", "Australia", "Beauty & Personal Car"] 37507429.66715513  
["2023-12-26T20", "Australia", "Clothing, Shoes & Jewelry"] 37505669.50929021  
["2023-12-26T20", "Australia", "Electronics"] 37507901.240210295  
["2023-12-26T20", "Australia", "Health, Household & Baby Care"] 37507723.693394884  
["2023-12-26T20", "Australia", "Home & Kitchen"] 37506749.883878775  
["2023-12-26T20", "Australia", "Office Supplies"] 37507907.86400904  
["2023-12-26T20", "Australia", "Pet Supplies"] 37505637.796687685  
["2023-12-26T20", "Australia", "Sports & outdoors"] 37506709.31300333  
["2023-12-26T20", "Australia", "Toys & games"] 37507649.93278636  
hadoop_user@hadoop-master:~$
```

Этап 3

- Установка и настройка Airflow
- Создание DAG для выполнения созданных ранее скриптов
- Выполнение DAG run несколько раз для доказательства работоспособности

Namenode information

All Applications

DAGs - Airflow

+

localhost:8080/home

Airflow

DAGs

Cluster Activity

Datasets

Security

Browse

Admin

Docs

10:08 UTC

UH

Do not use **SQLite** as metadata DB in production – it should only be used for dev/testing. We recommend using Postgres or MySQL. [Click here](#) for more information.

Do not use the **SequentialExecutor** in production. [Click here](#) for more information.

DAGs

All 1Active 0Paused 1Running 0Failed 0

Filter DAGs by tag

Search DAGs

Auto-refresh

<div></div> DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
<div><div></div><div>purchases_profit</div><div>business</div></div>	airflow	<div><div></div><div>8</div><div></div><div></div><div></div></div>	0:20:00	2023-12-30, 09:00:00	2023-12-31, 09:20:00	<div><div></div><div></div><div></div><div></div><div></div><div></div><div>4</div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div>

«

«

1

>

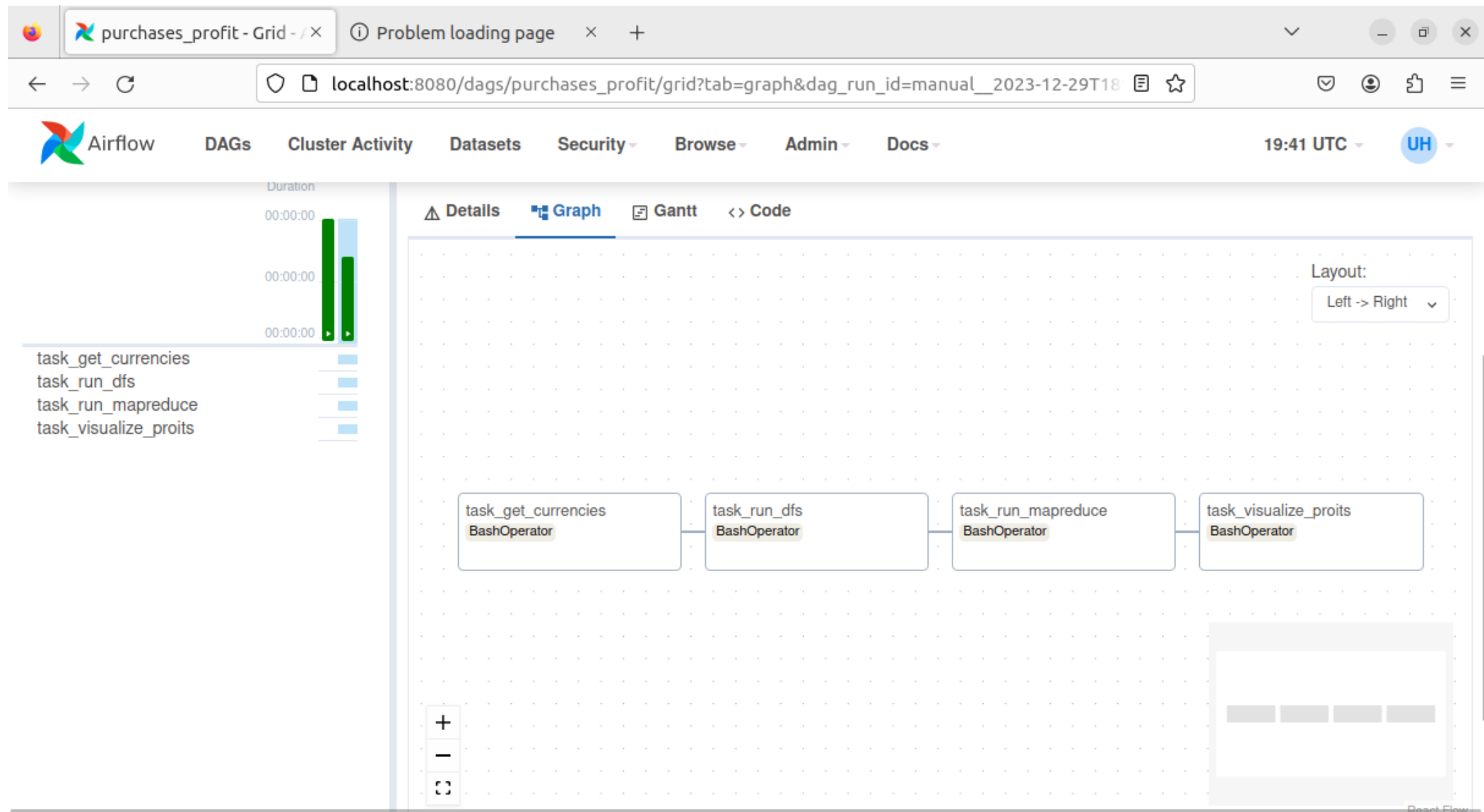
>>

Showing 1-1 of 1 DAGs

Version: v2.8.0

19

DAG



Namenode information

All Applications

purchases_profit - Grid - /x

localhost:8080/dags/purchases_profit/grid

10:10 UTC

UH

Airflow

DAGs

Cluster Activity

Datasets

Security

Browse

Admin

Docs

12 / 31 / 2023 , 10 : 10 : 06 AM

25

All Run Types

All Run States

Clear Filters

Auto-refresh

Press **shift** + **/** for Shortcuts

deferred

failed

queued

removed

restarting

running

scheduled

skipped

success

up_for_reschedule

up_for_retry

upstream_failed

no_status

task_get_currencies

task_run_dfs

task_run_mapreduce

task_visualize_proits

00:15:57

00:07:58

00:00:00

Dec 29, 23:40

DAG

purchases_profit

Details

Graph

Gantt

Code

DAG Runs Summary

Total Runs Displayed

9

Total success

8

Total running

1

First Run Start

2023-12-29, 22:01:39 UTC

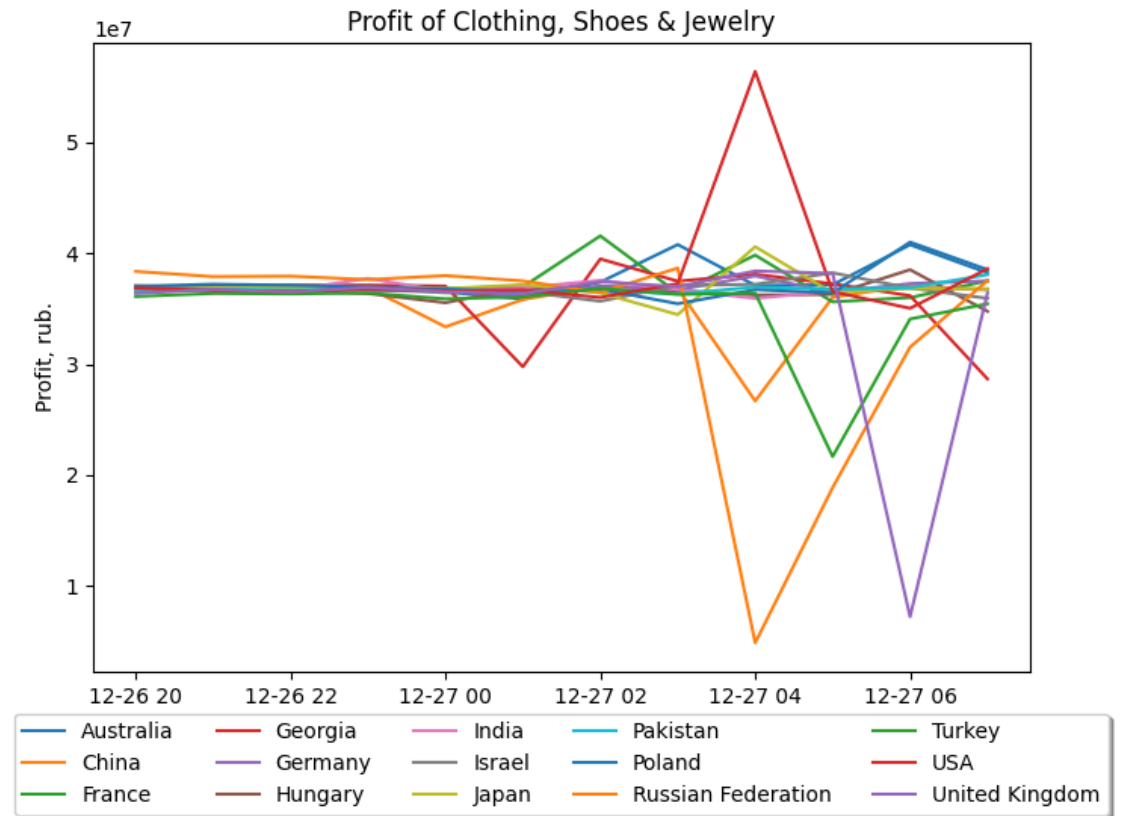
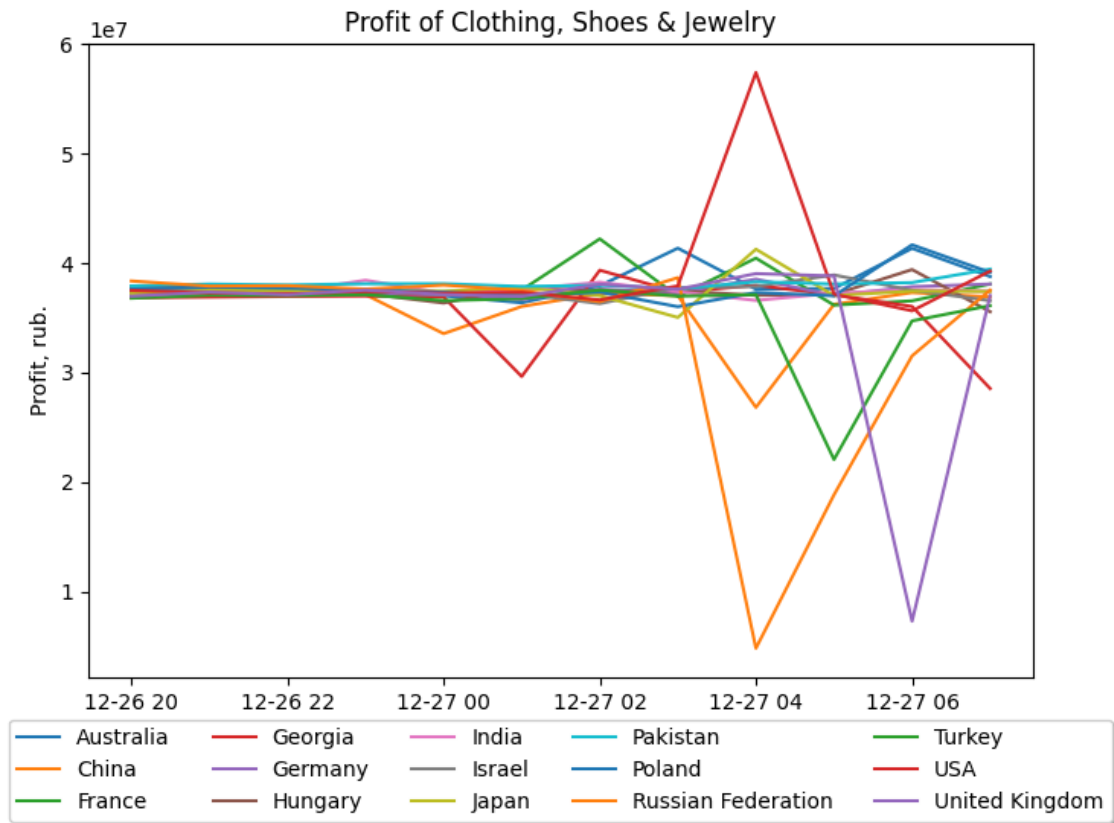
Last Run Start

2023-12-31, 10:08:09 UTC

Max Run Duration

00:15:57

График за 7 и 8 прогон



(разница есть)

Польза лабораторной

- Закрепление знаний по работа с терминалом, Bash, сетевыми службами, виртуализацией
- Навыки создания и настройки кластера
- Навыки оркестрации данных
- Навыки простого анализа данных
- Навыки работы с MapReduce

ИСТОЧНИКИ

1. Хабр. Big Data от А до Я. Часть 1: Принципы работы с большими данными, парадигма MapReduce [Электронный ресурс]. - URL: <https://habr.com/ru/articles/267361/> (дата обращения: 30.12.2023)
2. Tom White. Hadoop: The Definitive Guide [Электронный ресурс]. - URL: [ссылка](#) (дата обращения: 30.12.2023)
3. Apache Hadoop 3.3.6 Documentation [Электронный ресурс]. - URL: <https://hadoop.apache.org/docs/stable/> (дата обращения: 30.12.2023)
4. Apache Hadoop YARN Documentation [Электронный ресурс]. - URL: <https://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/YARN.html> (дата обращения: 30.12.2023)
5. Apache Airflow Documentation [Электронный ресурс]. - URL: <https://airflow.apache.org/docs/> (дата обращения: 30.12.2023)
6. Хабр. Все, что вам нужно знать об Airflow DAGs, ч.1 — Основы и расписания [Электронный ресурс]. - URL: <https://habr.com/ru/articles/682384/> (дата обращения: 30.12.2023)
7. mrjob v0.7.4 documentation [Электронный ресурс]. - URL: <https://mrjob.readthedocs.io/en/latest/> (дата обращения: 30.12.2023)
8. pandas 2.1.4 documentation [Электронный ресурс]. - URL: <https://pandas.pydata.org/docs/> (дата обращения: 30.12.2023)
9. Matplotlib 3.8.2 documentation [Электронный ресурс]. - URL: <https://matplotlib.org/stable/index.html> (дата обращения: 30.12.2023)
10. Mendel Cooper. Advanced Bash-Scripting Guide [Электронный ресурс]. - URL: https://se.ifmo.ru/~ad/Documentation/ABS_Guide_ru.html (дата обращения: 30.12.2023)

Результат работы

Лабораторная работа -

https://github.com/staffeev/cs_itmo_project

Видео-демонстрация работы кластера -

https://drive.google.com/file/d/16BC8VS1cVjWzKh3Dgtd5pA0_921YW1lD/view?usp=sharing

Спасибо за внимание!