Лабораторная работа 1. Установка и настройка распределенной системы. Простейшие операции и знакомство с функциональностью системы.

Цель: ознакомление с процессом установки и настройки распределенных систем, таких как Apache(Arenadata) Hadoop или Apache Spark. Изучить основные операции и функциональные возможности системы, что позволит понять принципы работы с данными и распределенными вычислениями.

Необходимое ПО:

- Ubuntu 24.04 LTS (22.04, 20.04) или новее.
- Java 8 ил Java11 или новее.
- Apache Spark 3.4.3.
- Python 3.12+.
- pip (менеджер пакетов Python).

Алгоритм выполнения задания в Apache Spark

- 1. Установка Java
- 2. Установка Python и рір
- 3. Установка Apache Spark 3.4.3
- 4. Настройка переменных окружения
- 5. Загрузка экономических данных
- 6. Запуск Spark и выполнение простейших операций

1. Установка Java.

```
"bash sudo apt update sudo apt install openjdk-11-jdk java -version
```

2. Установка Python и рір.

```
"bash sudo apt install python3 python3-pip python3 --version pip3 --version
```

3. Установка Apache Spark 3.4.3:

```
```bash
```

wget https://downloads.apache.org/spark/spark-3.4.3/spark-3.4.3-bin-hadoop3.tgz tar xvf spark-3.4.3-bin-hadoop3.tgz sudo mv spark-3.4.3-bin-hadoop3 /opt/spark

4. Настройка переменных окружения:

```
Откройте файл ~/.bashrc:
```

```
```bash
nano ~/.bashrc
```

Добавьте следующие строки в конец файла:

```
***Dash
export SPARK_HOME=/opt/spark
export PATH=$PATH:$SPARK_HOME/bin
```

• • •

Сохраните файл и примените изменения:

```
```bash
```

```
source ~/.bashrc
```

...

5. Загрузка данных.

Создайте директорию для данных и загрузите исторические данные по акциям (например, Apple):

```
```bash
```

mkdir ~/spark_data

cd ~/spark_data

wget -O AAPL.csv

или скачать по ссылке необходимый файл в **Downloads** и скопировать фвайл в необходимую директорию.

```bash

cp ~/Downloads/ AAPL.csv ~/spark_data

٠.,

6. Запуск Spark и выполнение простейших операций:

Запустите PySpark:

```bash

pyspark

Py ```

В интерактивной оболочке PySpark выполните следующие операции:

а. Загрузка данных.

```
```python
```

df = spark.read.csv("file:///home/username/spark_data/AAPL.csv", header=True,
inferSchema=True)

df.show(5)

...

```
>>> df = spark.read.csv("file:///home/devops/spark_data/AAPL.csv", header=True, inferSchema=True)
>>> df.show(5)

| Date| Open| High| Low| Close|Adj Close| Volume|
| 1980-12-12|0.128348|0.128906|0.128348|0.128348| 0.098943|469033600|
| 1980-12-15| 0.12221| 0.12221|0.121652|0.121652| 0.093781|175884800|
| 1980-12-16|0.113281|0.113281|0.112723|0.112723| 0.086898|105728000|
| 1980-12-17|0.115513|0.116071|0.115513|0.115513| 0.089049| 86441600|
| 1980-12-18|0.118862| 0.11942|0.118862|0.118862| 0.091631| 73449600|
```

```
b. Подсчет количества строк.
"python

print("Количество строк:", df.count())
""

c. Вывод схемы данных.
"python

df.printSchema()
""

>>> print("Количество строк:", df.count())
Количество строк: 11016
>>> df.printSchema()
root
|-- Date: date (nullable = true)
|-- High: double (nullable = true)
|-- Low: double (nullable = true)
|-- Adj Close: double (nullable = true)
|-- Volume: long (nullable = true)
```

d. Базовая статистика.

```python

df.describe().show()

• • •

е. Фильтрация данных.

```python

```
df_filtered = df.filter(df["Date"] >= "2020-01-01")
df_filtered show(F)
```

df_filtered.show(5)

``

```
f. Группировка и агрегация.
 ```python
 from pyspark.sql.functions import year, avq
 df_yearly = df.withColumn("Year",
year(df["Date"])).groupBy("Year").agg(avg("Close").alias("Avg_Close"))
 df_yearly.orderBy("Year").show()
```

```
1980|0.13590307692307693
1981|0.10854781818181822
```

g. Создание временного представления и выполнение SQL-запроса:

```
```python
```

df.createOrReplaceTempView("stock_data") spark.sgl("SELECT Year(Date) as Year, AVG(Close) as Avg Close FROM stock data GROUP BY Year(Date) ORDER BY Year").show()

++	
Year	Avg_Close
++	
1980 0.1359030	7692307693
1981 0.1085478	1818181822
1982 0.0854588	8142292486
1983 0.167274	0118577075
1984 0.119651	2490118577
1985 0.0902334	8809523808
1986 0.1449128	5375494065
1987 0.347751	1106719368
1988 0.370884	2411067196
1989 0.371952	9682539681
1990 0.3353729	0513834017
1991 0.4687014	7826086955
1992 0.489309	4763779526
1993 0.3663096	4426877505
1994 0.304290	5634920634
1995 0.3619659	5238095214
1996 0.2224953	4645669278
1997 0.160417	6679841897
1998 0.272900	4801587301
1999 0.51580	5492063492
++	+

Чтобы выйти из среды PySpark, выполните следующие действия:

1. Если вы работаете в интерактивной сессии PySpark (например, через терминал или Jupyter Notebook), просто введите:

```
```python
```

spark.stop()

Это завершит ceaнс Spark.

2. После этого можете выйти из Python интерактивной оболочки:

## ```python exit()

Или просто нажмите `Ctrl + D`. Это завершит сессию PySpark и закроет текущий Python интерпретатор.

### Алгоритм выполнения задания в Apache Hadoop

В виртуальной машине Шаг 1-8 пропустить.

Шаг 1. Установка необходимых компонентов.

```bash

```
sudo apt update
sudo apt install ssh pdsh -y
Шаг 2. Создание пользователя Hadoop.
sudo adduser hadoop
sudo usermod -aG sudo hadoop
su - hadoop
Шаг 3. Настройка SSH.
```bash
ssh-keygen -t rsa -P ""
cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
chmod 0600 ~/.ssh/authorized_keys
Шаг 4. Загрузка и установка Hadoop.
```bash
wget https://downloads.apache.org/hadoop/common/hadoop-3.3.5/hadoop-3.3.5.tar.qz
tar -xzvf hadoop-3.3.5.tar.gz
sudo mv hadoop-3.3.5 /usr/local/hadoop
Шаг 5. Настройка окружения Hadoop.
Добавьте следующие строки в конец файла ~/.bashrc:
```bash
export HADOOP_HOME=/usr/local/hadoop
export HADOOP INSTALL=$HADOOP HOME
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export HADOOP_COMMON_HOME=$HADOOP_HOME
export HADOOP HDFS HOME=$HADOOP HOME
export YARN_HOME=$HADOOP_HOME
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native
export PATH=$PATH:$HADOOP_HOME/sbin:$HADOOP_HOME/bin
Применить изменения:
```bash
source ~/.bashrc
Шаг 6. Настройка конфигурационных файлов Hadoop.
а) Отредактируйте $HADOOP_HOME/etc/hadoop/hadoop-env.sh:
```bash
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64
b) Отредактируйте $HADOOP HOME/etc/hadoop/core-site.xml:
```xml
```

```
<configuration>
  cproperty>
    <name>fs.defaultFS</name>
    <value>hdfs://localhost:9000</value>
  </property>
</configuration>
c) Отредактируйте $HADOOP_HOME/etc/hadoop/hdfs-site.xml:
```xml
<configuration>
 cproperty>
 <name>dfs.replication</name>
 <value>1</value>
 </property>
 cproperty>
 <name>dfs.namenode.name.dir</name>
 <value>/home/hadoop/hdfs/namenode</value>
 </property>
 property>
 <name>dfs.datanode.data.dir</name>
 <value>/home/hadoop/hdfs/datanode</value>
 </property>
</configuration>
d) Отредактируйте $HADOOP_HOME/etc/hadoop/mapred-site.xml:
```xml
<configuration>
  property>
    <name>mapreduce.framework.name</name>
    <value>yarn</value>
  </property>
</configuration>
e) Отредактируйте $HADOOP_HOME/etc/hadoop/yarn-site.xml:
```xml
<configuration>
 cproperty>
 <name>yarn.nodemanager.aux-services</name>
 <value>mapreduce_shuffle</value>
 </property>
</configuration>
Шаг 7. Создание директорий для HDFS.
```bash
mkdir -p ~/hdfs/namenode ~/hdfs/datanode
```

```
Шаг 8. Форматирование HDFS.
```bash
hdfs namenode -format
Все дальнейшие действия выполняются пользователем hadoop.
```bash
sudo su - hadoop
Шаг 9. Запуск Hadoop.
```bash
start-dfs.sh
start-yarn.sh
 hadoop@devopsvm:-$ jps
Шаг 10. Проверка работы Hadoop.
 5682 SecondaryNameNode
```bash
                                                 5940 ResourceManager
jps
                                                 6068 NodeManager
                                                 11236 SparkSubmit
                                                 5465 DataNode
                                                 9145 SparkSubmit
                                                 10250 SparkSubmit
                                                 5181 NameNode
                                                 12478 Jps
                                                 9422 SparkSubmit
```

Вы должны увидеть следующие процессы: NameNode, DataNode, SecondaryNameNode, ResourceManager, NodeManager.

В стандартной конфигурации Hadoop HDFS предоставляет веб-интерфейс, доступный через веб-браузер на порту 9870. Этот интерфейс позволяет просматривать состояние и структуру HDFS, а также выполнять некоторые операции.

Чтобы получить доступ к веб-интерфейсу HDFS, выполните следующие шаги:

- Убедитесь, что Hadoop (в частности, HDFS) запущен.
- Откройте веб-браузер на компьютере, с которого у вас есть сетевой доступ к серверу Hadoop.
- В адресной строке браузера введите:

http://localhost:9870



Overview 'localhost:9000' (ractive)

Started:	Mon Aug 26 13:07:51 +0300 2024
Version:	3.3.5, r706d88266abcee09ed78fbaa0ad5f74d818ab0e9
Compiled:	Wed Mar 15 18:56:00 +0300 2023 by stevel from branch-3.3.5
Cluster ID:	CID-60a52b68-6139-4947-8731-3c039547a32e
Block Pool ID:	BP-1830111676-127.0.1.1-1724666841903

- Если обращаетесь к Hadoop с другого компьютера, замените "localhost" на IP-адрес или имя хоста сервера, на котором запущен Hadoop.
- Нажмите Enter, должны увидеть веб-интерфейс HDFS.

Через этот веб-интерфейс вы сможете просматривать структуру директорий HDFS, проверять состояние и здоровье узлов, просматривать логи и выполнять другие административные задачи.

Шаг 11. Работа с экономическими данными

а) Создайте директорию в HDFS:

```bash

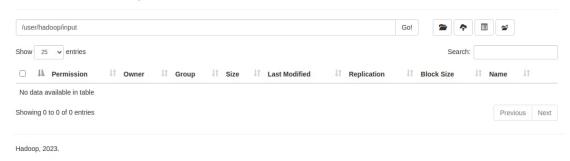
hdfs dfs -mkdir /user

hdfs dfs -mkdir /user/hadoop

hdfs dfs -mkdir /user/hadoop/input

` ` `

Browse Directory



Проведем расчет экономических показателей на примере открытых экономических данных. Будем использовать данные о ВВП стран мира от Всемирного банка.

Шаг 12. Подготовка данных.

Скачайте данные о ВВП стран мира

(https://raw.githubusercontent.com/BosenkoTM/Distributed_systems/main/practice/2024/lw_01/GDP.csv. Coxpаните файл как **GDP.csv**.

Шаг 13. Загрузка данных в HDFS.

```bash

wget

https://raw.githubusercontent.com/BosenkoTM/Distributed_systems/main/practice/2024/lw_01/G DP.csv

hdfs dfs -mkdir /user/hadoop/economic data

hdfs dfs -put GDP.csv /user/hadoop/economic_data/

Обработка данных с помощью MapReduce или Spark

Шаг 14. Запустите Spark.

```bash

spark-shell

Шаг 15. Загрузите данные и выполните расчеты.

```scala

val data = spark.read.option("header", "true").csv("file:///home/hadoop/GDP.csv")

Вычисляем среднее значение GDR

val result = data.selectExpr("avg(GDR) as avg_GDR")

```
val data = spark.read.option("header","true").csv("file:///home/hadoop/GDP.csv")
 lata: org.apache.spark.sql.DataFrame = [Country: string, Year: string ... 21 more fields]
scala> val result = data.selectExpr("avg(GDP) as avg GDR")
 result: org.apache.spark.sql.DataFrame = [avg_GDR: double]
// Сохраняем результат в CSV файл
result.write.option("header", "true").csv("/home/hadoop/output/avg_GDR.csv")
Выходим из Scala.
  ```scala
:q
 В Ubuntu 24.04 Scala сохраняет результаты в файле part-00000-*.csv, каталог будет
определен последним адресом в пути при сохранении, то есть avg GDR.csv.
Шаг 16. Переименовать полученный результат part-00000-*.csv в Ubuntu avg.csv
  ```bash
mv part-00000-*.csv avg.csv
Шаг 17. Переносим данные в HDFS. Загрузите экономические данные в HDFS:
```bash
hdfs dfs -put /home/hadoop/output/avg.csv /user/hadoop/input/
 Проверьте, что данные загружены.
```bash
hdfs dfs -ls /user/hadoop/input/
 /user/hadoop/input
                                                                 Go!
Show 25 v entries
 ☐ ♣ Permission
                11 Owner
                                  11 Size
                                       Last Modified
                                                      Replication
                                                                  ☐ Block Size
                                           Aug 26 15:42
                                                                                        â
                  hadoon
                                    27 B
                                                        1
                                                                    128 MB
      -LM-L--L--
                           supergroup
                                                                                avd.csv
       Чтобы остановить Hadoop 3 в Ubuntu, выполните следующие шаги.
Сначала остановите YARN (если он запущен).
```bash
stop-yarn.sh
Затем остановите HDFS.
```bash
stop-dfs.sh
       Если вы используете MapReduce JobHistory Server, остановите его:
       mr-jobhistory-daemon.sh stop historyserver.
       Для полной остановки всех Hadoop-демонов можно использовать команду.
```bash
stop-all.sh
Проверьте, что все процессы Hadoop остановлены.
```bash
ips
```

Эта команда покажет список запущенных Java-процессов. Убедитесь, что в списке нет процессов, связанных с Hadoop (например, NameNode, DataNode, ResourceManager и т.д.).

Если какие-то процессы остались, вы можете остановить их вручную с помощью команды kill.

```bash

kill -9 <PID>

где <PID> - идентификатор процесса, который вы хотите остановить.

#### Задание для самостоятельной работы

- 1. Загрузите данные по акциям другой компании (например, Microsoft MSFT).
- 2. Выполните аналогичный анализ для новых данных.
- 3. Сравните результаты анализа двух компаний.
- 4. Напишите Spark-приложение, которое находит дни с максимальным объемом торгов для обеих компаний.

#### Отчет

Студенты должны подготовить отчет, включающий:

- 1. Описание процесса установки и настройки Spark или Hadoop.
- 2. Листинг выполненных команд и их результаты.
- 3. Анализ полученных результатов.
- 4. Код и результаты выполнения задания для самостоятельной работы.
- 5. Выводы о функциональности и возможностях Apache Spark или Hadoop для анализа экономических данных.

#### Варианты заданий

Вариант выбирается согласно номеру студента в списке группы:

1. Установка Apache Spark на одном узле и выполнение простой задачи на подсчет строк в файле.

Данные: Исторические данные по акциям Сбербанка (SBER) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2020 год, расчет средней цены закрытия, группировка по месяцам.

2. Установка Apache(Arenadata) Наdoop и выполнение задачи на копирование файлов в HDFS.

Данные: Исторические данные по акциям Газпрома (GAZP) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2019 год, расчет максимальной цены открытия, группировка по кварталам.

3. Установка Apache Spark и выполнение задачи на сортировку данных.

Данные: Исторические данные по акциям Лукойла (LKOH) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 5 лет, расчет минимальной цены закрытия, группировка по годам.

4. Настройка кластерного режима для Apache Spark на 2 узлах.

Данные: Исторические данные по акциям Яндекса (YNDX) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2021 год, расчет средней цены закрытия, тренд анализа.

5. Настройка кластерного режима для Apache(Arenadata) Наdoop на 2 узлах и проверка работоспособности.

Данные: Исторические данные по акциям Роснефти (ROSN) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 3 года, расчет медианной цены закрытия, группировка по месяцам.

6. Установка Apache Spark на кластере из 3 узлов и выполнение задачи на агрегацию данных. Данные: Исторические данные по акциям Норильского никеля (GMKN) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2020 год, расчет стандартного отклонения цены закрытия, группировка по кварталам.

7. Установка и настройка Apache Spark для работы с внешним источником данных (например, S3).

Данные: Исторические данные по акциям ВТБ (VTBR) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за последние 10 лет, расчет коэффициента вариации цены закрытия, тренд

8. Установка Apache(Arenadata) Наdoop и выполнение задачи на объединение файлов в HDFS.

Данные: Исторические данные по акциям Магнита (MGNT) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2018 год, расчет средней цены открытия, группировка по годам.

9. Установка и настройка Apache Spark для работы с Cassandra.

Данные: Исторические данные по акциям Полюса (PLZL) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за 2019 год, расчет средней цены закрытия, корреляция с объемом торгов.

10. Настройка Apache Spark для работы с SQL-запросами.

Данные: Исторические данные по акциям MTC (MTSS) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за последние 2 года, расчет максимальной цены закрытия, тренд анализа.

11. Установка Apache(Arenadata) Наdoop и выполнение задачи на создание и удаление каталогов в HDFS.

Данные: Исторические данные по акциям Татнефти (TATN) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 5 лет, расчет минимальной цены закрытия, группировка по месяцам.

12. Установка Apache Spark и выполнение задачи на чтение и запись данных из/в Parquet.

Данные: Исторические данные по акциям Сургутнефтегаз (SNGS) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2020 год, расчет стандартного отклонения цены открытия, тренд анализа.

- 13. Настройка кластерного режима для Apache Spark на 4 узлах с разными ролями узлов. Данные: Исторические данные по акциям Мечела (MTLR) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за последние 3 года, расчет медианной цены закрытия, группировка по кварталам.
- 14. Установка Apache(Arenadata) Наdoop и выполнение задачи на распределение файлов между узлами.

Данные: Исторические данные по акциям Интер РАО (IRAO) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2018 год, расчет средней цены открытия, корреляция с объемом торгов.

15. Установка Apache Spark и выполнение задачи на анализ текстовых данных.

Данные: Исторические данные по акциям Аэрофлота (AFLT) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 2 года, расчет максимальной цены закрытия, тренд анализа.

16. Настройка кластерного режима для Apache Spark на 3 узлах с использованием Docker.

Данные: Исторические данные по акциям Системы (AFKS) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2019 год, расчет средней цены закрытия, группировка по месяцам.

17. Установка Apache(Arenadata) Наdoop и выполнение задачи на создание и просмотр логов системы.

Данные: Исторические данные по акциям ФосАгро (PHOR) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 5 лет, расчет минимальной цены закрытия, тренд анализа.

18. Установка Apache Spark и выполнение задачи на фильтрацию данных.

Данные: Исторические данные по акциям Алросы (ALRS) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за 2020 год, расчет стандартного отклонения цены закрытия, группировка

19. Настройка кластерного режима для Apache(Arenadata) Наdoop на 4 узлах и выполнение задачи на распределенную обработку данных.

Данные: Исторические данные по акциям Русала (RUAL) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за последние 3 года, расчет медианной цены закрытия, корреляция с объемом торгов

20. Установка Apache Spark и выполнение задачи на работу с JSON-файлами.

Данные: Исторические данные по акциям Мосбиржи (MOEX) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за 2018 год, расчет средней цены открытия, группировка по годам.

21. Настройка Apache Spark для работы с Hive.

Данные: Исторические данные по акциям РУСГИДРО (HYDR) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 10 лет, расчет коэффициента вариации цены закрытия, тренд анализа.

22. Установка Apache(Arenadata) Наdoop и выполнение задачи на создание резервной копии данных.

Данные: Исторические данные по акциям Россетей (RSTI) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за 2021 год, расчет максимальной цены закрытия, корреляция с объемом торгов.

23. Установка Apache Spark и выполнение задачи на вычисление статистических параметров данных.

Данные: Исторические данные по акциям X5 Retail Group (FIVE) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 2 года, расчет средней цены закрытия, группировка по месяцам.

- 24. Настройка кластерного режима для Apache Spark на 2 узлах с использованием Ansible. Данные: Исторические данные по акциям ТМК (TRMK) с сайта Московской биржи (moex.com) Операции: Фильтрация данных за 2019 год, расчет минимальной цены закрытия, тренд анализа.
- 25. Установка Apache(Arenadata) Наdoop и выполнение задачи на слияние данных из нескольких источников в HDFS.

Данные: Исторические данные по акциям М.Видео (MVID) с сайта Московской биржи (moex.com)

Операции: Фильтрация данных за последние 5 лет, расчет стандартного отклонения цены закрытия, группировка по кварталам.