Exercise:

- 1. Install spark and pyspark (ubuntu only)
- 2. Run a spark shell and test the installation.
- 3. Run the wordcount program that you did using hadoop using pyspark.
- 4. Use the movielens dataset available in the LMS theory page and try to find out for each movie, how are the ratings distributed.

Link:

https://medium.com/@agusmahari/pyspark-step-by-step-guide-to-installing-pyspark-on-linux-bb8af96ea5e8

1. Install spark and pyspark (ubuntu only)

- PySpark, a Python library for Apache Spark
- Below steps is to install Apache Spark and using pip install the PySpark library

Step 1: Install Java

java -version

```
shruthimohan@EDITH:~$ java -version
openjdk version "11.0.24" 2024-07-16
OpenJDK Runtime Environment (build 11.0.24+8-post-Ubuntu-1ubuntu320.04)
OpenJDK 64-Bit Server VM (build 11.0.24+8-post-Ubuntu-1ubuntu320.04, mixed mode, sharing)
```

Step 2. Install Apache Spark

wget

https://archive.apache.org/dist/spark/spark-3.2.0/spark-3.2.0-bin-hadoo
p3.2.tgz

extracted file: tar -xvzf spark-3.2.0-bin-hadoop3.2.tgz

```
### Cated file: tar -xvzf spark-3.2.0-bin-
3.2.0-bin-hadoop3.2/subernetes/
3.2.0-bin-hadoop3.2/subernetes/
3.2.0-bin-hadoop3.2/subernetes/
3.2.0-bin-hadoop3.2/subernetes/tests/
3.2.0-bin-hadoop3.2/subernetes/tests/
3.2.0-bin-hadoop3.2/subernetes/tests/
3.2.0-bin-hadoop3.2/subernetes/tests/subscale.py
3.2.0-bin-hadoop3.2/subernetes/tests/subscale.py
3.2.0-bin-hadoop3.2/subernetes/tests/subscale.py
3.2.0-bin-hadoop3.2/subernetes/tests/subscale.py
3.2.0-bin-hadoop3.2/subernetes/tests/decommissioning_cleanup.py
3.2.0-bin-hadoop3.2/subernetes/tests/decommissioning_cleanup.py
3.2.0-bin-hadoop3.2/subernetes/tests/decommissioning_cleanup.py
3.2.0-bin-hadoop3.2/subernetes/tests/decommissioning_cleanup.py
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/decom.sh
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/R/Docerfite
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/R/Docerfite
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/R/Docerfite
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/R/Docerfite
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/pomeraledecerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernetes/decerfites/spark/bindings/python/
3.2.0-bin-hadoop3.2/subernete
                                                                                                                                                                                                                                                                                                                                                                                                                                             CI-Desktop-PC:~/hadoop$ tar -xvzf spark-3.2.0-bin-hadoop3.2.tgz
```

Move the extracted folder to the /opt directory: sudo mv

```
spark-3.2.0-bin-hadoop3.2 /opt/spark
```

Step 3. Set Up Environment Variables

To set permanent environment variables for a single user klick on home in terminal

sudo nano ~/.bashrc

```
GNU nano 4.8
# this, if it's already enabled in /etc/bash.bashrc and /etc/profile
# sources /etc/bash.bashrc).
if ! shopt -oq posix; then
    if [ -f /usr/share/bash-completion/bash_completion ]; then
        . /usr/share/bash-completion/bash_completion
    elif [ -f /etc/bash_completion ]; then
        . /etc/bash_completion
fi

fi

export NVM_DIR="$HOME/.nvm"
[ -s "$NVM_DIR/nvm.sh" ] && \. "$NVM_DIR/nvm.sh" # This loads nvm
[ -s "$NVM_DIR/nvm.sh" ] && \. "$NVM_DIR/nvm.sh" # This loads nvm
export SPARK_HOME=/opt/spark
export PATH=$PATH:$SPARK_HOME/bin:$SPARK_HOME/sbin
```

```
hadoop@snucse-HP-Pro-Tower-400-G9-PCI-Desktop-PC:-/hadoop$ sudo nano -/.bashrc
hadoop@snucse-HP-Pro-Tower-400-G9-PCI-Desktop-PC:-/hadoop$ source -/.bashrc
hadoop@snucse-HP-Pro-Tower-400-G9-PCI-Desktop-PC:-/hadoop$ []
```

Step 4. Install PySpark

Install PySpark using pip: pip install pyspark

```
hadoop@snucse-HP-Pro-Tower-400-G9-PCI-Desktop-PC:-/hadoop$ pip Install pyspark
Defaulting to user installation because normal site-packages is not writeable
collecting pyspark
Downloading pyspark-3.5.3.tar.gz (317.3 MB)

Preparing metadata (setup.py) ... done
Collecting py4j=e0.10.9.7 (from pyspark)
Downloading py4j-e0.10.9.7 -py2.py3-none-any.whl.metadata (1.5 kB)
Downloading py4j-0.10.9.7 -py2.py3-none-any.whl (200 kB)

Downloading py4j-0.10.9.7 -py2.py3-none-any.whl (200 kB)

Building wheels for collected packages: pyspark
Building wheel for pyspark (setup.py) ... done
Created wheel for pyspark: filename-pyspark-3.5.3-py2.py3-none-any.whl size=317840630 sha256=ad620f0845e816c2194eab84cefa892c78d7517543bf64b213d4ddd2f29b33e6
Stored in directory: /home/hadoop/.cache/pip/wheels/97/f5/c0/947e2c0942b361ffe58651f36bd7f13772675b3863fd63d1b1
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully Installed py4j-0.10.9.7 pyspark-3.5.3
hadoop@snucse-HP-Pro-Tower-480-G9-PCI-Desktop-PC:-/hadoop$
```

Step 5. Verify PySpark Installation: pyspark

```
hadooppssucse.NP-Pro-Tower-400-60-PCL-Desktop-PC:-/hadoox; pyspark
Python 3.11.7 (main, Dec 15 2023, 1812:31) [SCC 11.20] on linux
Type "help", "copyright", "credits" or "license" for more information.
24/19/08 10:12:25 MARN Utils: Your hostname, snucse-HP-Pro-Tower-400-60-PCL-Desktop-PC resolves to a loopback address: 127.0.1.1; using 10.23.22.170 instead (on interface eno1)
24/19/08 10:12:25 MARN Utils: Set SPARK LOCAL_IP if you need to bind to another address
Using Spark's default logi-porfile: org/apache/spark/log4j-defaults.properties
Setting default log level to "MARN".
To adjust logging level use sc.settoglevel(newLevel). For SparkR, use settoglevel(newLevel).
24/19/08 10:12:25 MARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Welcome to

Version 3.1.7 (main, Dec 15 2023 18:12:31)
Spark context Web UI available at http://10.23.22.170:4040
Spark context Web UI available at http://10.23.22.170:4046
Spark scontext Web UI available as 'sc' (master = local[*], app id = local-1728362546168).
>>> 

SparkSession available as 'spark'.
```

2. Run a spark shell and test the installation.

Run pyspark:

```
hadoop@EDITH:~/word_count$ pyspark

Python 3.9.5 (default, Nov 23 2021, 15:27:38)

[GCC 9.3.0] on linux

Type "help", "copyright", "credits" or "license" for more information.

24/10/10 10:49:11 WARN Utils: Your hostname, EDITH resolves to a loopback address: 127.0.1.1; using 172.22.79.226 instead (on interfa ce etho)

24/10/10 10:49:11 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address

WARNING: An illegal reflective access operation has occurred

WARNING: An illegal reflective access by org.apache.spark.unsafe.Platform (file:/opt/spark/jars/spark-unsafe_2.12-3.2.0.jar) to construc tor java.nio.DirectByteBuffer(long,int)

WARNING: Please consider reporting this to the maintainers of org.apache.spark.unsafe.Platform

WARNING: See --illegal-access=warn to enable warnings of further illegal reflective access operations

WARNING: All illegal access operations will be denied in a future release

Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties

Setting default log level use sc.setlogLevel(newLevel). For SparkR, use setLogLevel(newLevel).

24/10/10 10:49:12 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where a pplicable

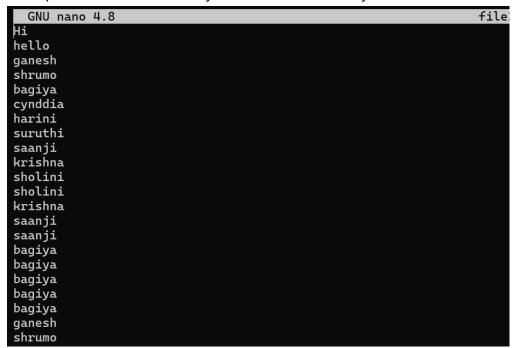
Welcome to
```

3. Run the wordcount program that you did using hadoop using pyspark.

Link: https://www.learntospark.com/2020/01/word-count-program-in-apache-spark.html

Step 1: cd word_count (directory created in prev ex)

The input file 'file1.txt' is already created in this directory.



Step 2: create a SparkSession and sparkContext

Step 3: Read the input file as RDD and provide transformations to calculate the count of each word in our file.

Step 4: Output

NOTES:

- 1. SparkSession is the entry point for creating a Spark application.
- 2. .master("local"): Specifies that the Spark application should run in "local" mode, meaning it will use local threads instead of connecting to a Spark cluster
- 3. sc = spark.sparkContext: Gets the Spark context (sc) from the SparkSession. sc is the main entry point for Spark's RDD-based API, providing access to core Spark functions like data loading, transformations, and actions.
- 4. RDD= Resilient Distributed Dataset
- 5. flatMap(lambda line: line.split(" ")): Applies a function to each line in text_file, splitting it into individual words. flatMap flattens the result so each word becomes a separate element in the RDD.

Program done in jupyter notebook:

```
[7]: import pyspark
      # Create SparkSession and sparkcontext
      from pyspark.sql import SparkSession
      spark = SparkSession.builder\
                          .master("local")\
                          .appName('wordcount')\
                           .getOrCreate()
      sc=spark.sparkContext
 [8]: # Read the input file and Calculating words count
      text_file = sc.textFile("word_count/file1.txt")
      counts = text_file.flatMap(lambda line: line.split(" ")) \
                                   .map(lambda word: (word, 1)) \
                                  .reduceByKey(lambda x, y: x + y)
[10]: output = counts.collect()
      for (word, count) in output:
          print('{}\t{}'.format(word, count))
      Ηi
      hello 1
      ganesh 2
      shrumo 2
      bagiya 9
      cynddia 4
      harini 4
      suruthi 4
      saanji 3
      krishna 2
      sholini 4
[11]: # Stopping Spark-Session and Spark context
      sc.stop()
      spark.stop()
```

4. Use the movielens dataset available and try to find out for each movie, how are the ratings distributed.

Code:

from pyspark.sql import SparkSession

from pyspark.sql.functions import col, count

Initialize Spark session

```
spark = SparkSession.builder \
    .appName("MovieLensRatingsDistribution") \
    .getOrCreate()

# Load dataset

data_path = "movielens.txt"

movies_df = spark.read.csv(data_path, sep='\t', inferSchema=True) \
    .toDF("user_id", "movie_id", "rating", "timestamp")

# Calculate the distribution of ratings for each movie

rating_distribution = movies_df.groupBy("movie_id", "rating") \
    .agg(count("rating").alias("rating_count")) \
    .orderBy("movie_id", "rating")
```

Show the results

rating_distribution.show()

```
# Show the results
rating distribution.show()
+----+
|movie_id|rating|rating_count|
+----+
   1 1 8
               27
    1 2
    1 3
               96
    1 4
              202
    1 5
               119
    2
        1
               8
    2
        2
               17
    2
        3
                55
    2
        4
               42
    2
        5
                9
    3
        1
                11
    3
        2
                20
    3
        3
                25
    3
        4
                23
    3
        5
               11
    4
        1
                6
    4
        2
                24
    4
        3
                57 l
    4
        4
                931
        5
    4
+----+
only showing top 20 rows
```