

Computer Science & Information Systems

Big Data Systems – Spark Lab Sheet 5

Spark MLlib

1. Objective:

Students should be able to

- A. Get familiarity with the Spark MLlib module
- B. Get hands-on experience with Machine Learning programme development and execution

This lab sheet provides a quick introduction of using Spark for Machine Learning applications. This exercise will introduce the API through MLlib package for development of classical regression model.

MLlib is Spark's machine learning (ML) library. Its goal is to make practical machine learning scalable and easy. At a high level, it provides tools such as:

- ML Algorithms: common learning algorithms such as classification, regression, clustering, and collaborative filtering
- Featurization: feature extraction, transformation, dimensionality reduction, and selection
- Pipelines: tools for constructing, evaluating, and tuning ML Pipelines
- Persistence: saving and load algorithms, models, and Pipelines
- Utilities: linear algebra, statistics, data handling, etc.

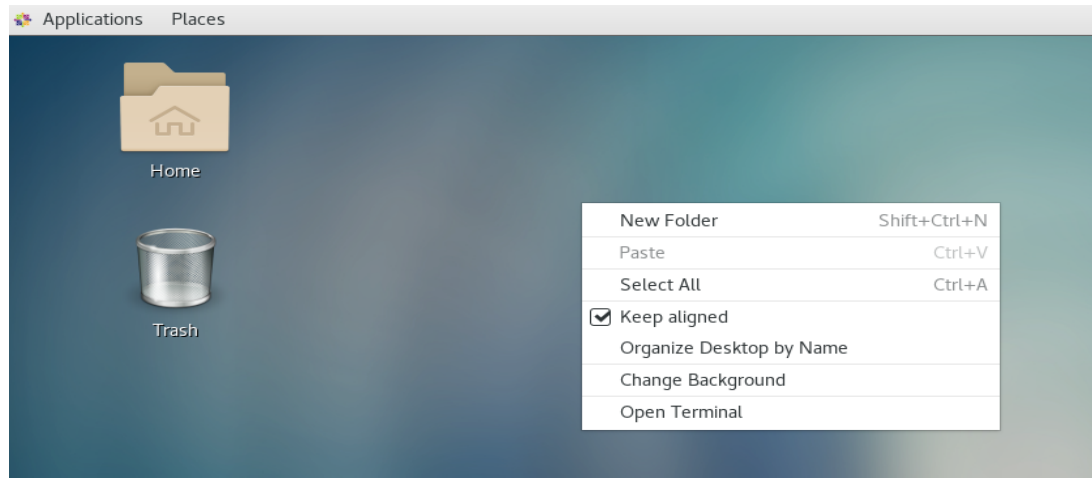
2. Steps to be performed:

Note - It's assumed that student has made a slot reservation using the slot booking interface where Apache Spark framework was selected. The details of the Apache Spark systems to be used is received through an email. If not, please contact the administrators for the same.

Also it's assumed that students are aware of the process of logging into these virtual machines. If not, then get access to the user manual maintained for the usage of remote lab setup.

Preparations -

- a) Open the terminal by right clicking on the desktop of the virtual machine.



- b) Look at the current directory and also file listings in it. It must have a spark installation directory present in it. Commands like pwd, ls can be used for it.

```

Applications  Places  Terminal
csishydlab@apache-spark:~

File Edit View Search Terminal Help
[csishydlab@apache-spark ~]$ pwd
/home/csishydlab
[csishydlab@apache-spark ~]$ ls
boston.csv          k1.py               people.json          spark-2.4.4-bin-hadoop2.7
consumer.py          kafka_2.12-2.4.0    Pictures              spark-2.4.4-bin-hadoop2.7.tgz
d1.py                kafka_2.12-2.4.0.tgz producer.py           Spark-DataFrame.ipynb
data                  log.txt              Public                spark-streaming-kafka-0-8-assembly_2.11-2.4.4.jar
Desktop              lr1.py               sample_linear_regression_data.txt
direct_kafka_wordcount.py Music                 sample_svm_data.txt  spark-warehouse
Documents             network_wordcount.py scala-2.10.1.tgz      students.json
Downloads             output               sk1.py               Templates
input.txt             p1.py               Videos
[csishydlab@apache-spark ~]$

```

- c) Set the SPARK_HOME and HOME variable to point to the spark installations.

```
[csishydlab@apache-spark bin]$ pwd
```

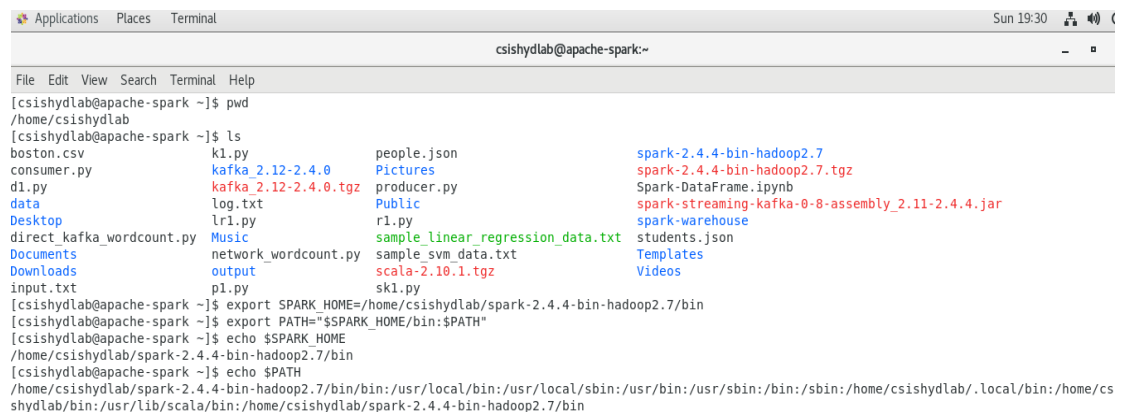
```
/home/csishydlab/spark-2.4.4-bin-hadoop2.7/bin
```

```
[csishydlab@apache-spark bin]$ export SPARK_HOME=/home/csishydlab/spark-2.4.4-bin-hadoop2.7/bin
```

```
[csishydlab@apache-spark bin]$ export PATH="$SPARK_HOME/bin:$PATH"
```

```
echo $SPARK_HOME
```

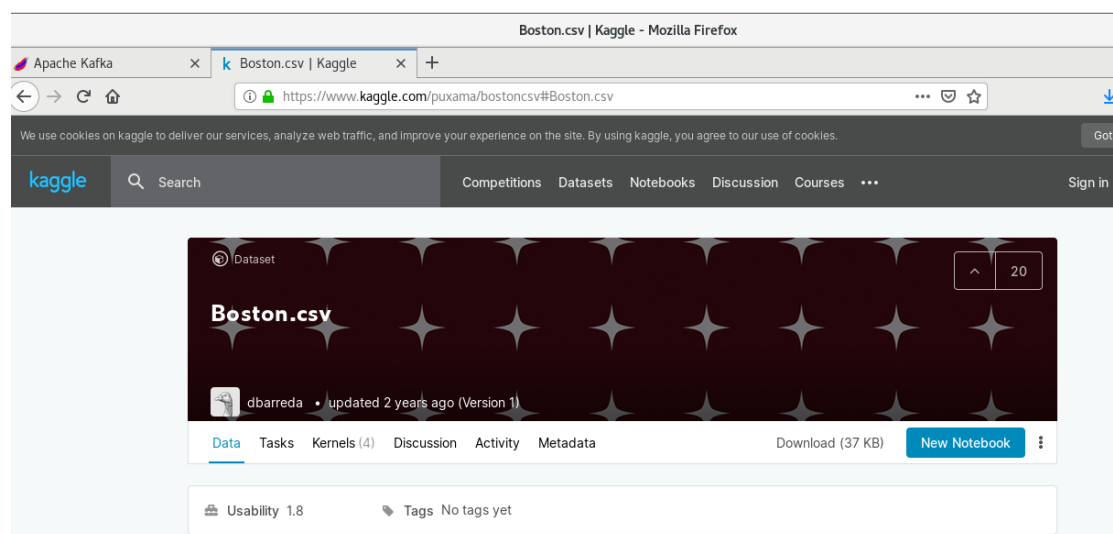
```
echo $PATH
```



```
csishydlab@apache-spark:~$ pwd
/home/csishydlab
csishydlab@apache-spark:~$ ls
boston.csv          kafka_2.12-2.4.0    people.json          spark-2.4.4-bin-hadoop2.7
consumer.py         kafka_2.12-2.4.0.tgz Pictures              spark-2.4.4-bin-hadoop2.7.tgz
d1.py               log.txt             producer.py           Spark-DataFrame.ipynb
data                lr1.py              Public                spark-streaming-kafka-0-8-assembly_2.11-2.4.4.jar
Desktop             network_wordcount.py sample_linear_regression_data.txt spark-warehouse
direct_kafka_wordcount.py Music               sample_svm_data.txt  students.json
Documents           output              scala-2.10.1.tgz     Templates
Downloads           pl.py              sk1.py               Videos
input.txt           network_wordcount.py sample_svm_data.txt  students.json
input.txt           output              scala-2.10.1.tgz     Templates
input.txt           pl.py              sk1.py               Videos
[csishydlab@apache-spark ~]$ export SPARK_HOME=/home/csishydlab/spark-2.4.4-bin-hadoop2.7/bin
[csishydlab@apache-spark ~]$ export PATH="$SPARK_HOME/bin:$PATH"
[csishydlab@apache-spark ~]$ echo $SPARK_HOME
/home/csishydlab/spark-2.4.4-bin-hadoop2.7/bin
[csishydlab@apache-spark ~]$ echo $PATH
/home/csishydlab/spark-2.4.4-bin-hadoop2.7/bin:/usr/local/bin:/usr/local/sbin:/usr/bin:/usr/sbin:/bin:/sbin:/home/csishydlab/.local/bin:/home/csishydlab/bin:/usr/lib/scala/bin:/home/csishydlab/spark-2.4.4-bin-hadoop2.7/bin
```

- d) Using the web browser , download the Boston.csv file and save it in the local file system.

<https://www.kaggle.com/puxama/bostoncsv#Boston.csv>



Installing pySpark

For the execution of python programmes on the Spark, a package named pyspark is required. Using the sudo privileges, install the packages with pip command.

```
pip install pyspark
```

Writing Linear Regression Machine Learning programme

- e) Open up the text editor and copy the code written in the attached linear_regression.py file.

```
csishydlab@
File Edit View Search Terminal Help
[root@apache-spark csishydlab]# gedit linear_regression.py &
```

- f) Execute the linear_regression.py file using the spark-submit command.

```
csishydlab@
File Edit View Search Terminal Help
[root@apache-spark csishydlab]# gedit linear_regression.py &
[3] 25797
[2] Done gedit linear_regression.py
[root@apache-spark csishydlab]# spark-submit linear_regression.py
```

- g) Look at the outcome printed while the program is getting executed on the Spark cluster. It shows actual and predicted house prices.

```
+-----+-----+-----+
|          prediction | medv |          features |
+-----+-----+-----+
| 27.927270908006726 | 22.0 | [0.01096,55.0,2.2... |
| 30.768510069596182 | 29.1 | [0.01439,60.0,2.9... |
| 26.497762769897413 | 23.1 | [0.0187,85.0,4.15... |
| 40.29308710051661 | 50.0 | [0.02009,95.0,2.6... |
| 27.087495861316143 | 16.5 | [0.02498,0.0,1.89... |
| 28.147245589564633 | 23.9 | [0.02543,55.0,3.7... |
| 22.219573905760136 | 20.6 | [0.03306,0.0,5.19... |
| 32.57075672757935 | 34.9 | [0.03359,75.0,2.9... |
| 20.383617525016447 | 19.5 | [0.03427,0.0,5.19... |
| 30.645810541006668 | 28.5 | [0.03502,80.0,4.9... |
| 27.895366964633475 | 22.0 | [0.03537,34.0,6.0... |
| 24.551339355753125 | 22.9 | [0.03551,25.0,4.8... |
| 28.756201579966294 | 27.9 | [0.03615,80.0,4.9... |
| 23.549490273145608 | 20.7 | [0.03738,0.0,5.19... |
| 35.0073271643155 | 34.6 | [0.03768,80.0,1.5... |
| 36.411263074979686 | 33.3 | [0.04011,80.0,1.5... |
| 26.88971216116669 | 22.9 | [0.04203,28.0,15... |
| 24.69531323293066 | 20.6 | [0.04294,28.0,15... |
| 26.554206975416243 | 24.8 | [0.04297,52.5,5.3... |
| 17.14499460979746 | 18.2 | [0.04301,80.0,1.9... |
+-----+-----+-----+
only showing top 20 rows
```

Installing numpy

- h) For the execution of python programmes on the Spark, a package named numpy is required. Using the sudo privileges, install the packages with pip command.

```
pip install numpy
```

Writing k-means clustering Machine Learning programme

- i) Open up the text editor and copy the data written in the attached kmeans_data.txt file.

```
Applications  Places  Text Editor

Open  [icon]  kmeans_data.txt
/home/csishydlab

kmeans_clustering.py  x

0.0 0.0 0.0
0.1 0.1 0.1
0.2 0.2 0.2
9.0 9.0 9.0
9.1 9.1 9.1
9.2 9.2 9.2
```

- j) Open up the text editor and copy the code written in the attached kmeans_clustering.py file.

```
kmeans_clustering.py  x

"""
A K-means clustering program using MLlib.

This example requires NumPy (http://www.numpy.org/).
"""
from __future__ import print_function

import sys

import numpy as np
from pyspark import SparkContext
from pyspark.mllib.clustering import KMeans

def parseVector(line):
    return np.array([float(x) for x in line.split(' ')])

sc = SparkContext(appName="KMeans")
lines = sc.textFile("kmeans_data.txt")
data = lines.map(parseVector)
k = 2
model = KMeans.train(data, k)
print("*****Final centers: " + str(model.clusterCenters))
print("*****Total Cost: " + str(model.computeCost(data)))
sc.stop()
```

- k) Execute the kmeans-clustering.py file using the spark-submit command.

- l) Look at the outcome printed while the program is getting executed on the Spark cluster. It shows the centroids based on the number of clusters mentioned in the python code (k value) .

```
20/01/26 22:40:33 INFO ContextCleaner: Cleaned accumulator 120
20/01/26 22:40:33 INFO MapPartitionsRDD: Removing RDD 3 from persistence list
20/01/26 22:40:33 INFO BlockManagerInfo: Removed broadcast 8 piece0 on apache-spark:34724 in memory (size: 5.5
*****Final centers: [array([0.1, 0.1, 0.1]), array([9.1, 9.1, 9.1])]
20/01/26 22:40:33 INFO BlockManager: Removing RDD 3
```

- m) Repeat the execution by changing the value of k in the python code.

3. Outputs/Results:

Students should be able to

- Execute the linear regression programme on Spark cluster
- See the accuracy of the linear regression models in terms of error
- See the predictions done by the model for the test records
- Execute the k-means programme on Spark cluster
- Look at the centroid generated by the programme with different values of k

4. Observations:

Students carefully needs to observe

- Error statistics associated with the linear regression model
- Centroids change when the number of clusters changed



5. References:

- A. [Spark ML Guide](#)
- B. [Spark Documentation](#)
- C. [Linear Regression](#)
- D. [K-means clustering](#)