# INSTALLING SPARK ON COLAB

## DR. ALVIN ANG

# CONTENTS

- You may reference here for help:
  https://tatwan.github.io/blog/colab/python/spark/2020/01/06/Colab-Spark-Instructions.html

- IPYNB: https://www.alvinang.sg/s/How_To_Start_A_Spark_Session.ipynb

- But bear in mind that Apache Spark regularly changes its download link (as well as Spark upgrades), thus the code needs to be modified from time to time.

```
[19] !apt-get install openjdk-8-jdk-headless -qq > /dev/null
```

### A. CODE

# !apt-get install openjdk-8-jdk-headless -qq > /dev/null
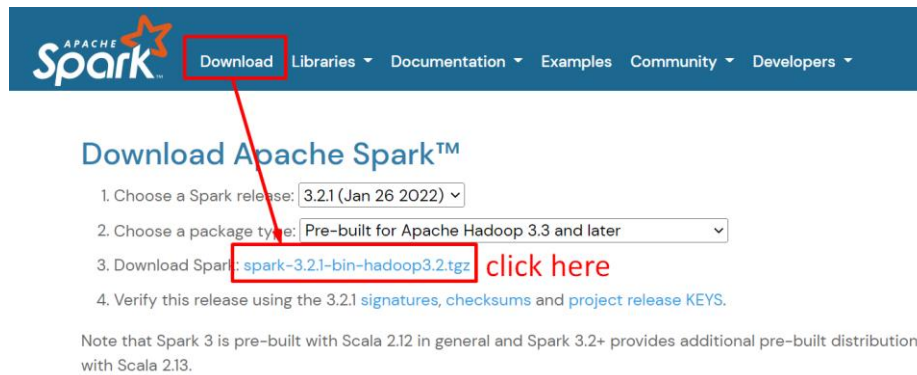
- Run and make sure this line of code works first.

```
[20] !wget -q https://dlcdn.apache.org/spark/spark-3.2.1/spark-3.2.1-bin-hadoop3.2.tgz
```

### A. CODE

**!wget -q https://dlcdn.apache.org/spark/spark-3.2.1/spark-3.2.1-bin-hadoop3.2.tgz**

- If this line of code doesn't work, most probably Apache has changed its link.

- Go here: https://spark.apache.org/downloads.html

```
✓  [22]  !tar xf spark-3.2.1-bin-hadoop3.2.tgz
4s
```

A. CODE

# !tar xf spark-3.2.1-bin-hadoop3.2.tgz

- Remember: If Apache upgrades its veresion to 3.X.X or whatever new versions, do edit the name accordingly.

```
✓ [23] !pip install -q findspark
3s
```

## A. CODE

# !pip install -q findspark

- This is to initialize PySpark.

- Findspark ensures that the environment variables will be properly set.

- PySpark will be imported upon importing findspark

```
[27] import os
     os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
     os.environ["SPARK_HOME"] = "/content/spark-3.2.1-bin-hadoop3.2"
```

## A. CODE

**import os**

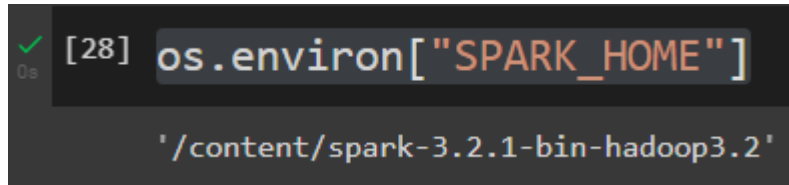**os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"**

**os.environ["SPARK_HOME"] = "/content/spark-3.2.1-bin-hadoop3.2"**

- Once again, remember to change the folder names to the latest updated Apache Spark version.

- Make sure every line runs properly before going to the next line of code.

## A. CODE

# os.environ["SPARK_HOME"]

- Just to check where the folder is…

Which can also be written as….



**there are many ways to write the Sparksession codes.. which we will see below…

# import findspark

# findspark.init()

# from pyspark.sql import SparkSession

# spark = SparkSession.builder.master("local[*]").getOrCreate()

*note:*

- *Config and Appname are non-essential... it can be ignored for now first...*

- *But those above are essential...*

Another way is to write the code below…

```
import findspark
findspark.init()  //search for SPARK and set it in the system path
from pyspark.sql import SparkSession //kickstart Spark Context
spark = SparkSession.builder.master("local[*]").getOrCreate() //since this is not a distributed cluster,
spark.conf.set("spark.executor.memory", "4g")                  driver and executor nodes are all local in Colab
spark.conf.set("spark.driver.memory", "4g")
spark.conf.set("spark.memory.fraction","0.9")  //setting up your Master and Slave memories.. if u want to...
```
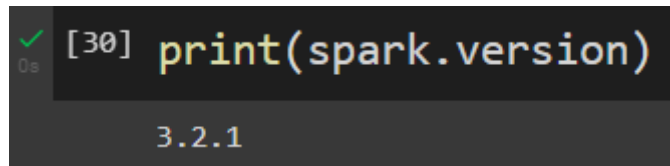
- findspark.init() = kickstart Spark Context

- SparkSession = to connect to Spark cluster

- Builder = to create a new SparkSession.

- Master =

    o Sets the Spark master URL to connect to, such as "local" to run locally, "local[4]" to run locally with 4 cores, or "spark://master:7077" to run on a Spark standalone cluster.

    o Usually, it would be either yarn or mesos depends on your cluster setup.

    o Use local[x] when running in Standalone mode (which is obvious in Colab since Colab can't run clusters..)

    o x should be an integer value and should be greater than 0.

    o Ideally, x value should be the number of CPU cores you have.

    o Putting as * is to let the computer find a value automatically for you, which is by default.

- getOrCreate = Gets an existing SparkSession or, if there is no existing one, creates a new one.

● If no configuration is done for the CPU cores / memory, it will be automatically set to MAX OUT your computer capacity…

## A.  CODE

# print(spark.version)

https://www.alvinang.sg/s/How_To_Start_A_Spark_Session.ipynb

### A. START A SPARK SESSION

```
[4]  !apt-get install openjdk-8-jdk-headless -qq > /dev/null

 ●   !wget -q https://dlcdn.apache.org/spark/spark-3.2.1/spark-3.2.1-bin-hadoop3.2.tgz

[6]  !tar xf spark-3.2.1-bin-hadoop3.2.tgz

[7]  !pip install -q findspark

[8]      import os
         os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"
         os.environ["SPARK_HOME"] = "/content/spark-3.2.1-bin-hadoop3.2"

[9]      os.environ["SPARK_HOME"]

     '/content/spark-3.2.1-bin-hadoop3.2'

[10]     import findspark
         findspark.init()

[11]     from pyspark.sql import SparkSession
         spark = SparkSession.builder.master("local[*]").getOrCreate()

[12]     print(spark.version)

     3.2.1
```

```
▾ Importing Libraries

✓ [13]  from pyspark.sql import SparkSession
   Os   spark=SparkSession.builder.appName('data_processing').getOrCreate()

        import pyspark.sql.functions as F
        from pyspark.sql.types import *
```

- appName =

  o Used to set your application name.

  o If no name is set, a randomly generated name will be used.

- Config = Sets a config option by specifying a (key, value) pair (used to setup your cores and memory I believe...)

Dr. Alvin Ang earned his Ph.D., Masters and Bachelor degrees from NTU, Singapore. He is a scientist, entrepreneur, as well as a personal/business advisor. More about him at www.AlvinAng.sg.