# Spark

> ℹ️ A Spark cheat sheet tailored to HOGWARTS can be found here.

**The following sections describe some of the most common SPARK use cases.**

> ⌄ Getting started - My first SPARK application
>
> ## To get started with SPARK, create a Notebook in jupyhub and follow these steps.
>
> 1. Create a basic SPARK context
>
> ```
> 1  from pyspark import SparkContext
> 2  from pyspark.sql import SparkSession
> 3  import os
> 4
> 5  spark = SparkSession.builder.getOrCreate()
> ```
>
> 2. Create a SPARK dataframe from an array of data and display its schema
>
> ```
> 1  arrayData = [
> 2          ('James',['Java','Scala'],{'hair':'black','eye':'brown'}),
> 3          ('Michael',['Spark','Java',None],{'hair':'brown','eye':None}),
> 4          ('Robert',['CSharp',''],{'hair':'red','eye':''}),
> 5          ('Washington',None,None),
> 6          ('Jefferson',['1','2'],{})
> 7  ]
> 8
> 9  df = spark.createDataFrame(data=arrayData, schema = ['name','knownLanguages','properties'])
> 10 df.printSchema()
> ```
>
> 3. Show the dataframe using vertical layout and no truncation
>
> ```
> 1  df.show(vertical=True, truncate=False)
> ```
>
> 4. Write the dataframe on disk (in /tmp/people_dir) in parquet format
>
> ```
> 1  %rm -rf /tmp/people_dir
> 2  df.write.parquet('file:///tmp/people_dir')
> 3  %ls -lh /tmp/people_dir
> ```
>
> 5. Read the data from disk into a second dataframe and show it
>
> ```
> 1  df2 = spark.read.parquet('file:///tmp/people_dir')
> 2  df2.show()
> ```
>
> 6. Stop the SPARK context
>
> ```
> 1  spark.stop()
> ```
>
> Congratulations, you have just completed your first SPARK application! Keep reading if you'd like to know how to access the datalake and/or the Hive data sources.

> ⌄ Using Spark UI in JupyterHub
>
> To start a Spark master run this:

```
1  $SPARK_HOME/sbin/start-master.sh
```

This will start a Spark master in your jupyterhub, it runs on port 8080, if you want to use a different port, you can add the option `--webui-port=<other port>`. This web ui is accessible in jupyterhub at the url:

```
1  https://jupyhub.hogwarts.<environment>.azure.chimera.cyber.gc.ca/user/<last_name>_<first_name>_<initial>_/pro
```

You may now be noticing that there's no worker to run your Spark jobs. Run this command on jupyterhub:

```
1  $SPARK_HOME/sbin/start-slave.sh spark://$HOSTNAME:7077
```

To use this Spark instance in your jobs, set the Spark master URI to the resolved version of `spark://$HOSTNAME:7077`. You can do this in a few ways:

1.
```
1  #!/bin/bash
2  echo "spark://$HOSTNAME:7077"
```

   Copy the output of this and paste it into your Spark master URI

2. Or you can do it all in python:

```
1  master_uri = f'spark://{os.environ.get("HOSTNAME")}:7077'
```

⌄ Accessing data sources stored in Hive

## To access the data sources stored as Hive tables, create a Notebook in jupyhub and follow these steps.

1. Create a SPARK context with support for HIVE

```
1  from pyspark import SparkContext
2  from pyspark.sql import SparkSession
3  from pyspark.sql.functions import *
4
5  spark = ( SparkSession.builder
6              .enableHiveSupport()
7              .getOrCreate()
8          )
```

2. List all tables that start with 'geo'

```
1  spark.sql("show tables like 'geo*'").show(truncate=False)
```

3. Print the schema of the geo_quova table

```
1  spark.sql("select * from geo_quova").printSchema()
```

4. Perform a SELECT statement on that table

```
1  spark.sql("select start_ip_int, end_ip_int, country from geo_quova").show(10)
```

5. Register scala User Defined Functions (UDFs) and use them to convert IPs to Strings

```
1  from hogwarts.spark.sql.functions import ipv4_long_col_to_str, register
2  register(spark)
3
4  spark.sql("select ipv4_long_col_to_str(start_ip_int), ipv4_long_col_to_str(end_ip_int), country from geo_quov
```

6. Stop SPARK context

```
1  spark.stop()
```

## Other useful pyspark sql methods

Rename columns using **as**:

```
1  spark.sql("select ipv4_long_col_to_str(start_ip_int) as start_ip, ipv4_long_col_to_str(end_ip_int) as end_ip,
```

Filter the results:

```
1  df = spark.sql("select ipv4_long_col_to_str(start_ip_int) as start_ip, ipv4_long_col_to_str(end_ip_int) as en
2  df.filter("country != 'united states'").show()
```

Convert a reasonable amount of rows (25 in the example below) to a Panda dataframe:

```
1  spark.sql("select ipv4_long_col_to_str(start_ip_int) as start_ip, ipv4_long_col_to_str(end_ip_int) as end_ip,
```

⚠️ IMPORTANT: trying to convert a huge spark dataframe to a Panda dataframe will result in out of memory error.

Describe Hive table including detailed table information:

```
1  spark.sql("describe formatted geo_quova").show(100, truncate=False)
```

˅ Reading from/Writing to the datalake

## To access the datalake using SPARK, create a Notebook in jupyhub and follow these steps.

### 1. Create a basic SPARK context

```
1  from pyspark import SparkContext
2  from pyspark.sql import SparkSession
3
4  spark = SparkSession.builder.getOrCreate()
```

### 2. Create a SPARK dataframe from an array of data and display its schema

```
1   arrayData = [
2          ("23.10.80.250", "e9631.j.akamaiedge.net"),
3          ("40.90.22.185", "fe-by01p-msa.com"),
4          ("104.16.121.127", "medium.com"),
5          ("69.16.175.10", "cds.r3t6j3w4.hwcdn.net"),
6          ("13.226.139.75", "ui.playboyengineering.com"),
7          ("31.13.71.49", "mmx-ds.cdn.whatsapp.net"),
8          ("104.16.151.254", "www.vibe.com")
9   ]
10
11  df = spark.createDataFrame(data=arrayData, schema = ['ip','hostname'])
12  df.printSchema()
```

### 3. Write a CSV to the **users** datalake container in the **demo/ip_to_hostname** folder

```
1  df.coalesce(1).write.mode("overwrite").csv("abfss://users@hogwartsdatalakeusersu.dfs.core.windows.net/demo/ip
```

### 4. Read back the CSV into a second dataframe called **df2** and show it

```
1  df2 = spark.read.format("csv").load("abfss://users@hogwartsdatalakeusersu.dfs.core.windows.net/demo/ip_to_hos
2  df2.printSchema()
```

```
3  df2.show(truncate=False)
```

5. Stop SPARK context

```
1  spark.stop()
```

## The SPARK cluster

In the previous sections, we have used the mini-spark cluster bundled within jupyhub. This is very useful to develop and test your spark application locally but this instance is limited in number of cores and memory. So when you need to scale your analytic to handle more data, you can launch your application on the HOGWARTS SPARK cluster. That cluster is equipped with dozens of cores and hundreds of GiB of memory. To run your application on the cluster, simply add the SPARK master uri when creating your SPARK context:

```
1  spark = ( SparkSession.builder
2             .master("spark://ver-1-spark-master-svc.spark:7077")
3             .appName("my awesome app")
4             .config("spark.executor.memory", "1g")
5             .config("spark.cores.max", '2')
6             .enableHiveSupport()
7             .getOrCreate()
8         )
```

Please note that you can specify the name of your SPARK application using `.appName`. This is the name that will appear in the SPARK console. The SPARK console can be accessed using a me2sa browser at: https://spark.u.hogwarts.azure.chi/ (https://spark.hogwarts.pb.azure.chimera.cyber.gc.ca/ on PB). You can look at the progress of your execution from that console.

> ⚠ IMPORTANT: when running your application on the SPARK cluster, you have to be careful about the number of resources that you request for your application (in order to share resources with your colleagues). You can limit the number of cores and memory that your application will use using the following config parameters: `"spark.executor.memory"` and `"spark.cores.max"`.

## Using custom environment for Spark Workers

If you wish to use a custom environment for specific UDF, there is a possibility to send a packed environment to the workers before executing your functions.

As found in the Spark documentation, starting with Spark 3.1 you can add the following two lines to your code to specify which file to use.

```
1  os.environ['PYSPARK_PYTHON'] = "./mycustomenv/bin/python"
2
3  spark = ( SparkSession.builder
4             .master("spark://ver-1-spark-master-svc.spark:7077")
5             .appName("my awesome app")
6             .config("spark.executor.memory", "1g")
7             .config("spark.cores.max", '2')
8             .config("spark.archives", f"abfss://users@{datalakename}.dfs.core.windows.net/path/to/custom_conda_
9             .enableHiveSupport()
10            .getOrCreate()
11        )
```

The name "mycustomenv" can be modified, but needs to be the same between the PYSPARK_PYTHON value and the value specified after the '#' in the spark.archives configuration.

You can use libraries like dlbrowse to upload files to Azure Datalake.

A quick reminder on Conda environment:

```
1   #You can create an environment using
2   conda create -y -n custom_conda_env
3   #NOTE: hogwarts-pyspark is currently using python 3.8, so you may need to force conda to create an environme
4   conda create -y -n custom_conda_env python=3.8
5   #You can activate your environment using
6   conda activate custom_conda_env
7   #You can add libraries using
8   conda install requests==2.22.0
9   conda install conda-pack
10  #After having conda-pack installed (either in your env or elsewhere), you can package it
11  conda pack -f -o custom_conda_env.tar.gz
12  #You can get out of a conda environment using
13  conda deactivate
```

If you want to use a different version of python, you need to use the same version for both the Spark driver (your notebook, or airflow) and the Spark worker.

```
1   #Create an environment specifying a different version of python
2   conda create -y -n python36_conda_env_worker python=3.6 conda-pack
3   conda activate python36_conda_env_worker
4   conda install requests==2.23.0 #Could be done at create time
5   conda pack -f -o python36_conda_env_worker.tar.gz
6   conda deactivate
7
8   #Create a second environment for your driver, which will contain libraries that are not needed by the worker
9   conda create -y -n python36_conda_env_driver python=3.6
10  conda activate python36_conda_env_driver
11  #Install a version of pyspark that match the one deployed by DASI2B
12  conda install pyspark
13  conda install ipykernel
14  #Install your new conda environment as a kernel for JupyterHub
15  python -m ipykernel install --user --name=python36_conda_env_driver
16
17  #You can now check that the kernel was added using
18  jupyter kernelspec list
19
20  #Then you can change your Notebook kernel to use python36_conda_env_driver
```

∨  Advanced topics

Create or replace a temporary view from a spark dataframe:

```
1   # Let's assume that we have a spark dataframe called df
2   df.createOrReplaceTempView('MY_TEMP_TABLE')
3   df2 = spark.sql('select count(*) from MY_TEMP_TABLE').show()
```

Create a User Defined Function (UDF):

```
1   df = spark.sql("select ip_long_to_str(start_ip_int) as start_ip, ip_long_to_str(end_ip_int) as end_ip, count
2
3   @udf("boolean")
4   def isInNorthAmerica(country):
5       if(country == "united states" or country == "canada" or country == "mexico"):
6           return True;
7       else:
8           return False;
9
10  df.withColumn("is_in_north_america", isInNorthAmerica(df.country)).show()
```

## References

- PySpark SQL documentation
- PySpark SQL Cheat Sheet
- PySpark Style Guide
- HOGWARTS built-in User Defined Functions
  - Python examples that use these UDFs
- HOGWARTS Spark Cheat Sheet
- Hogwarts PySpark API