

How to run Spark in Jupyter notebook on laptop and midway login and compute nodes

Igor Yakushin
ivy2@uchicago.edu

Laptop

- Assuming that you have installed latest java 1.8.x, Anaconda3, spark 2.4.4 on your laptop and they are in your path, to run spark in jupyter notebook do in a terminal

```
PYSPARK_DRIVER_PYTHON=jupyter PYSPARK_DRIVER_PYTHON_OPTS="notebook" pyspark
```

then in jupyter select New -> Python 3.

- On Mac and Linux you can check if you are using correct versions of java, python, spark by using `which`

```
which python
```

```
which java
```

```
which pyspark
```

and observe what paths are returned. For example, on my laptop the above commands return:

```
/usr/local/Anaconda3-2019.07/bin/python
```

```
/usr/local/jdk1.8.0_212/bin/java
```

```
/usr/local/spark-2.4.4-bin-hadoop2.7/bin/pyspark
```

- I believe, the corresponding command on Windows is `where`

Remote Linux system

- Suppose you want to run jupyter on midway or some other remote Linux system.
- Once you set your environment, for example by loading spark module, you can do the same command as on your laptop provided that your ssh client does X-forwarding.
- However, it would be slow because the web browser would be running on the remote system and the corresponding graphics would be transmitted via network.
- A better way would be to start a web server on the remote system and then connect the browser running on your laptop to the corresponding URL. The command from the previous slide is modified as follows:

```
PYSPARK_DRIVER_PYTHON=jupyter PYSPARK_DRIVER_PYTHON_OPTS="notebook --no-browser --ip=<host>" pyspark
```

where **<host>** is the hostname or ip address of the remote system.
- The above command would return you a URL to which you point the browser running on your laptop.

midway login nodes

- The command from the previous slide can be used to run spark in a jupyter notebook on midway login node.
- Make sure to provide the correct hostname.
- There are two login nodes on midway:
`midway2-login1.rcc.uchicago.edu` and
`midway2-login2.rcc.uchicago.edu`
- When you ssh to `midway2.rcc.uchicago.edu`, the load balancer puts you on one of them.
- However, it is not recommended to run anything heavy on login nodes. They exist for preparing and submitting batch jobs. They are shared by several thousand users and it is important not to overload them with computations.
- There is a script running on login nodes that would kill all the process of the heaviest user once the load on a node exceeds certain threshold.
- The best way to run spark on midway is to use compute nodes. There are about 500 of them on midway.

midway compute nodes

- To run jupyter notebook on a compute node, you need to use VPN or be on campus, since compute nodes are only visible on UChicago network and you cannot connect to the URL returned by jupyter running on the compute node from the outside of UChicago network.
- ssh to login node: `ssh <username>@midway2.rcc.uchicago.edu`
- Ask the scheduler for resources. For example, if you want to use all CPU cores (28) and all the memory (64G) in one node of broadwl partition for 3 hours:

```
sinteractive -p broadwl --exclusive --time=03:00:00 --nodes=1
```

- Note, if the node is available, you might get it within couple minutes. If it is not available, sinteractive would hang until there is an available node
- To check how busy a particular partition is:

```
sinfo -p broadwl
```

Look if there are any **idle** nodes.

midway compute nodes

- When you get a node, you are not sharing it with anybody. However, once you run out of the requested time, you are kicked out of the compute node even if the whole cluster is idle. So make sure to save your work periodically and ask for enough time. However, the more resources you ask for, the more time you might spend waiting in queue.
- Once you are on the compute node, you need to figure out its IP address since node name is not known outside of RCC cluster. The node's IP address is in the form `10.x.x.x`. To find it, use `hostname -i`.
- After that execute the same command as on login node but use IP address instead of hostname.

- To make life easier, I have written a script that you can use both on login and compute nodes. Execute on the node where you want to run spark in jupyter

```
source /project2/msca/ivy2/software2/etc/setup.sh
```

- The script loads Anaconda3, java 1.8, spark 2.4.4 and sets couple aliases.
- When you want to run spark in jupyter notebook on login node, execute `jnl`.
- When you want to run spark in jupyter notebook on compute node, execute `jncs`
- To see what each alias does, execute either `alias jnl` on login node or `alias jncs` on compute node.