# PySpark Test

Jungwon Seo

# Open Anaconda navigator and Launch Jupiter lab

# If you haven't installed pyspark, you can install using "conda install"

```
(base) ➜ Downloads conda install pyspark
Collecting package metadata: done
Solving environment: done

## Package Plan ##

  environment location: /Users/seojungwon/anaconda3

  added / updated specs:
    - pyspark


The following packages will be downloaded:

    package                    |            build
    ---------------------------|-----------------
    ca-certificates-2019.1.23  |                0         126 KB
    certifi-2019.3.9           |           py36_0         155 KB
    conda-4.6.8                |           py36_0         1.7 MB
    openssl-1.1.1b             |       h1de35cc_1         3.4 MB
    py4j-0.10.7                |           py36_0         250 KB
    pyspark-2.4.0              |           py36_0       203.5 MB
    ---------------------------|-----------------
                                        Total:       209.1 MB

The following NEW packages will be INSTALLED:

  py4j                 pkgs/main/osx-64::py4j-0.10.7-py36_0
  pyspark              pkgs/main/osx-64::pyspark-2.4.0-py36_0
```

# Or from Anaconda navigator

# You will see the familiar interface.
# Let's open terminal.

# Follow the GitHub instruction.

# Just remember to give the absolute path

# Check the result

# Good luck!