# How to use Hadoop & Spark?

# Contents

1. **Using Windows**

   - Anaconda installation

   - Spark installation

   - Hadoop installation

2. **Using Linux**

   - Virtual box installation
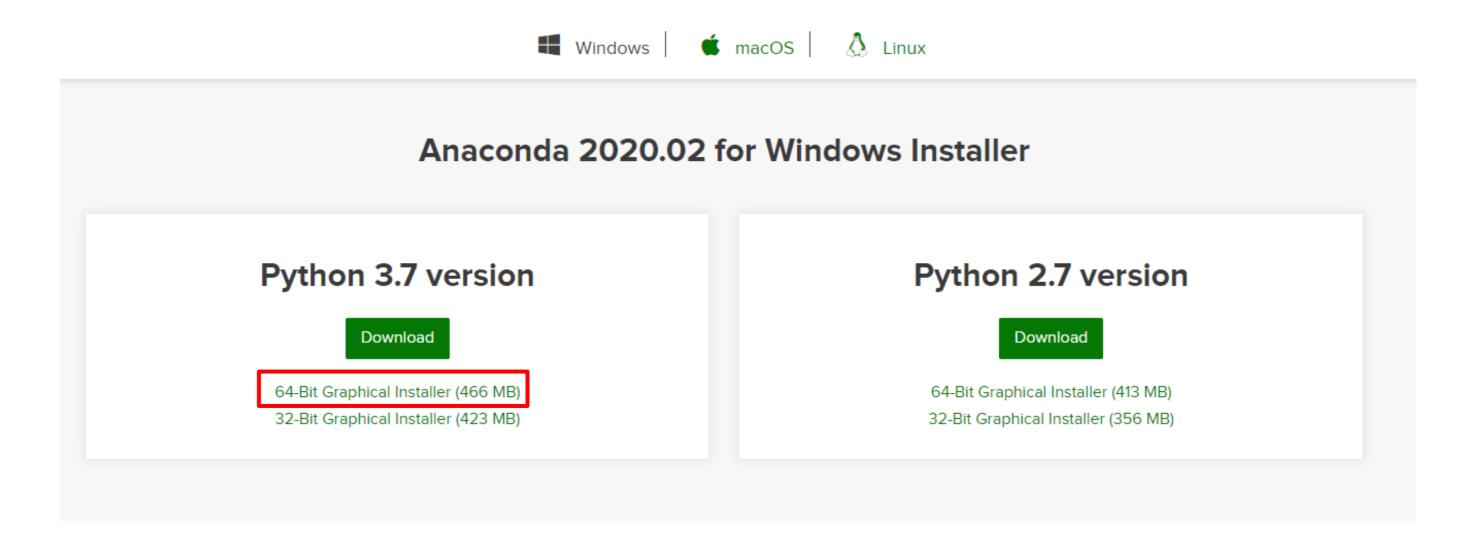
   - Spark installation

   - Hadoop installation

3. **Spark implementation**
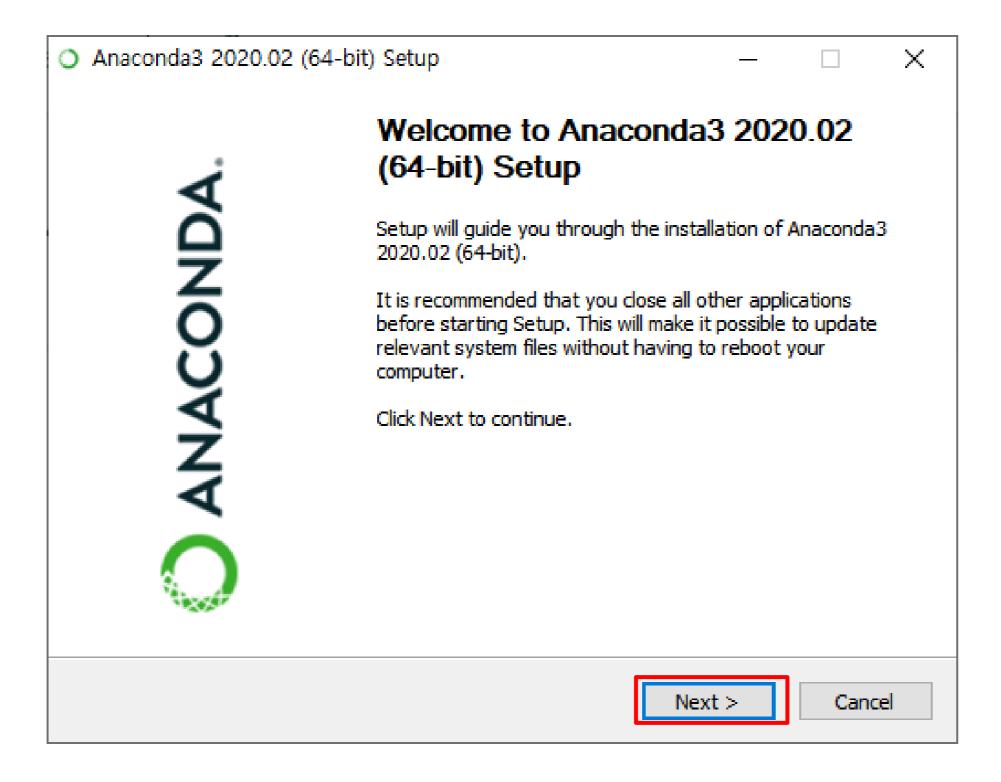
4. **Issues**

# Using Windows

# Anaconda Installation

# Install Anaconda

➢ **Anaconda is a tool for Python programming and command prompt**

Windows | macOS | Linux

### Anaconda 2020.02 for Windows Installer

**Python 3.7 version**

Download

64-Bit Graphical Installer (466 MB)
32-Bit Graphical Installer (423 MB)

**Python 2.7 version**

Download

64-Bit Graphical Installer (413 MB)
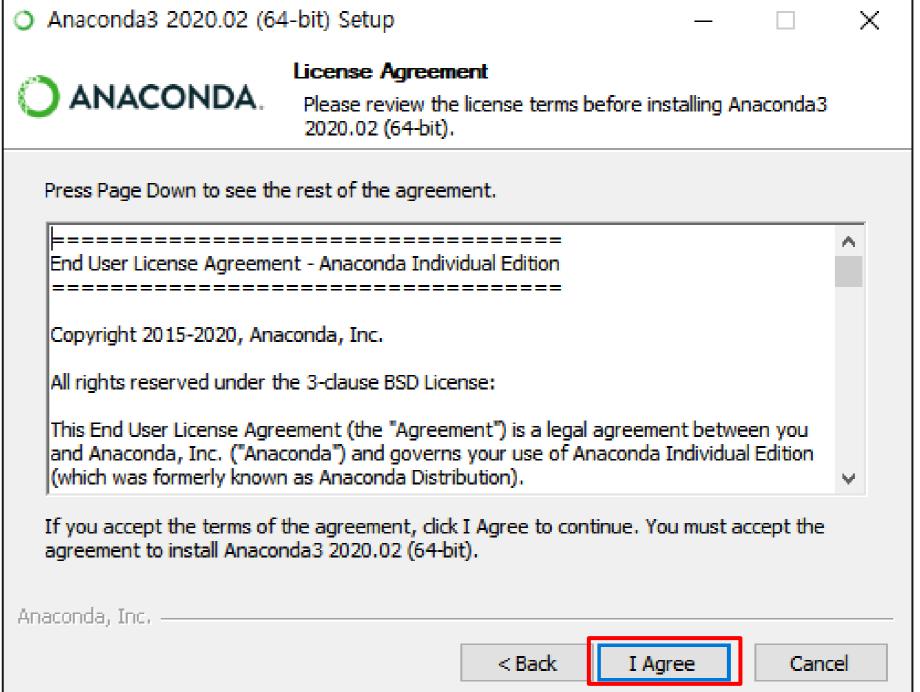32-Bit Graphical Installer (356 MB)
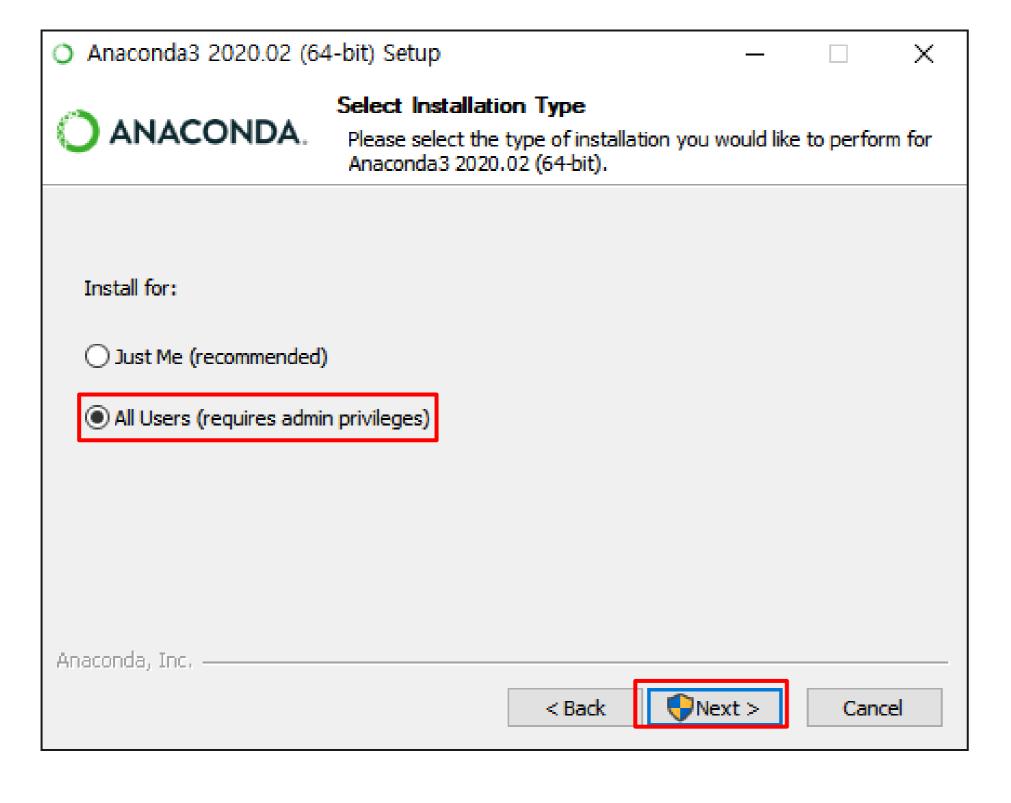
**Get Started with Anaconda Individual Edition**

➢ **If your computer's OS is Windows x86-64, We recommend download Python 3.7 64-bit Graphical Installer.**

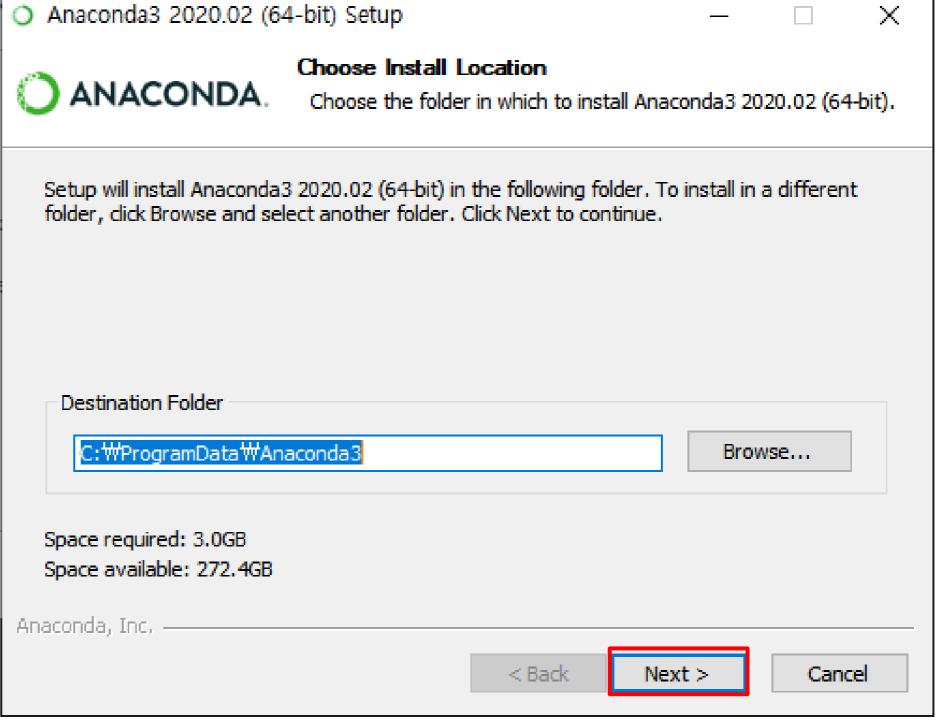➢ **You can download another version if your computer has another OS.**

https://www.anaconda.com/distribution/

5

# Install Anaconda

➢ **Begin installment and accept the license**

# Install Anaconda

➢ **Check "All Users" to avoid system authority issue**
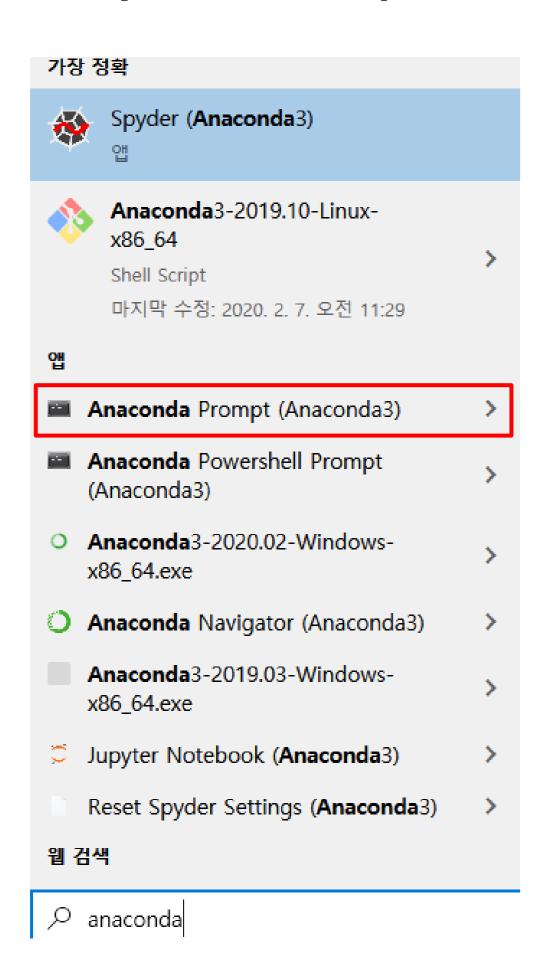
# Install Anaconda

➢ **Check "Register Anaconda3 as the system Python 3.7"**

➢ **Be careful not to check "Add Anaconda3 to the system PATH environment variable."**

➢ **Finish installment**

# Anaconda

- ➢ **Search Anaconda Prompt and open it**
- ➢ **We will use Python and Spark using this tool**

# Spark Installation

# Download Spark

➢ **Download Spark package from web site**

# Download Spark

➢ **Choose spark 2.4.5 & Apache Hadoop 2.7**



➢ **Then, download "spark-2.4.5-bin-hadoop2.7.tgz"**

➢ **Or, you can download the file here:**

https://www.apache.org/dyn/closer.lua/spark/spark-2.4.5/spark-2.4.5-bin-hadoop2.7.tgz

# Download Spark

➢ **Go to the folder where the package is installed and unzip the package for installment**

```
Anaconda Prompt (Anaconda3)

(base) C:\Users\     >tar zxvf spark-2.4.5-bin-hadoop2.7.tgz
```

➢ **After then, you can enter the "spark-2.4.5-bin-Hadoop.2.7" folder**

```
Anaconda Prompt (Anaconda3)

(base) C:\Users\     >cd spark-2.4.5-bin-hadoop2.7
(base) C:\Users\     \spark-2.4.5-bin-hadoop2.7>
```

# Download Spark

➤ **Type "bin\pyspark" and open Pyspark**

# Download Spark

➢ **Type "bin\spark-shell" and open Spark-Shell**

➢ **Spark-Shell is based on Scala**

# Configure System Path

➢ **We can configure system path for calling Spark in any folder.**

➢ **Enter Control Panel > System and Security > System**

# Configure System Path

➢ **System > Advanced System Settings**

# Configure System Path

➢ **"Advanced" tab > Environment Variables**

➢ **Add Spark folder to Path variable**



➢ **In our case, we add the path "C:\Users\BigdataLab\spark-2.4.5-bin-hadoop2.7\bin**

# Configure System Path

➤ **Return to Anaconda Prompt, and configure conda path with the following command**

**conda-develop SPARK_PATH\\python\\**

**conda-develop SPARK_PATH\\python\\lib\\**

```
(base) C:\Users\wjlee>conda-develop C:\Users\wjlee\spark-2.4.5-bin-hadoop2.7\python\
```
For example

```
(base) C:\Users\wjlee>conda-develop C:\Users\wjlee\spark-2.4.5-bin-hadoop2.7\python\lib\
```

➤ **Go to SPARK_PATH\python\lib, then drag py4j folder out from py4j-0.10.7-scr.zip**

# Configure System Path

➢ **Then, you can import pyspark library in Python**

```
Python 3.7.6 (default, Jan  8 2020, 20:23:39) [MSC v.1916 64 bit (AMD64)] :: Anaconda, Inc. on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import pyspark
>>>
```

# Hadoop installation

# Download JAVA

➢ **Download JAVA installation file (64bit)**

   *Refer the link on the bottom.*

Java를 다운로드하면 귀하가 <u>Oracle Java SE</u>에 대한 <u>Oracle Technology Network</u> 라이센스 합의서를 읽
고 이 조항에 동의하는 것으로 간주됩니다.

| | | | |
|---|---|---|---|
| Windows | ℹ <u>무엇을 선택해야 합니까?</u> | | |
| 🔴 <u>Windows 온라인</u><br>파일 크기: 1.97 MB | <u>지침</u> | Java를 설치한 후 브라우 |
| 🔴 <u>Windows 오프라인</u><br>파일 크기: 65.3 MB | <u>지침</u> | 저에서 Java를 사용으로<br>설정하려면 브라우저를<br>재시작해야 할 수 있습니 |
| 🔴 <u>Windows 오프라인 (64비트)</u><br>파일 크기: 73.29 MB | <u>지침</u> | 다. |

32비트 및 64비트 브라우저를 교대로 사용하는 경우, 각 브라우저에 대해 Java Plug-in이 필요하므로 32
비트 Java와 64비트 Java를 모두 설치해야 합니다. » <u>Windows용 64비트 Java에 대한 FAQ</u>

➢ **In our case, we install JAVA compatible with Windows OS.**

➢ **If your desktop has another OS, you need to install another version of JAVA.**

https://www.java.com/ko/download/manual.jsp

# Install Hadoop

➢ **Download pre-compiled package for Windows**

Pre-compiled, unofficial Win64 Binaries for Hadoop 2.7.1

| ⦿ 6 commits | ⑂ 1 branch | 📦 0 packages | 🏷 1 release | 👥 1 contributor | ⚖ Apache-2.0 |
|---|---|---|---|---|---|

| Branch: master ▾ | New pull request | | Create new file | Upload files | Find file | Clone or download ▾ |
|---|---|---|---|---|---|---|

| 👤 **karthikj1** Updated README | | Latest commit 910e032 on 10 Aug 2015 |
|---|---|---|
| 📄 Apache_README.txt | Added README | 5 years ago |
| 📄 LICENSE.txt | First commit | 5 years ago |
| 📄 NOTICE.txt | First commit | 5 years ago |
| 📄 README.md | Updated README | 5 years ago |
| 📄 build.png | First commit | 5 years ago |

📖 **README.md**

## Apache Hadoop 2.7.1 binary for Windows 64-bit platform

This is an unofficial pre-compiled binary of Apache Hadoop 2.7.1 for Windows 64-bit platform. The tar.gz file is available here under the Releases link in this repo.

The official Hadoop release from Apache does not include a Windows binary and compiling from sources can be tedious so I've made this compiled distribution available.

*https://github.com/karthikj1/Hadoop-2.7.1-Windows-64-binaries*

# Install Hadoop

➢ **In your Spark folder, make a new folder named Hadoop**



➢ **Open the Hadoop tar.zip file and move all folder or files into Hadoop folder**

*NOTE: We recommend there must be no space in the path of new folder*

*For example "C:\Users\big data class\spark\Hadoop folder" (X)*

# Configure PATH

➢ **Make new folder "C:\Hadoop" and Return to Anaconda prompt**

➢ **Link java jre file to "C:\Hadoop" with the following command**

*NOTE: There must be no SPACE in your path!*

*mklink /j C:\Hadoop\Java "C:\Program Files\Java\jre1.8.0_241"*

➢ **Then you can see system successfully link two paths like following**

```
(base) C:\Users\BigDataLab>mklink /j C:\Hadoop\Java "C:\Program Files\Java\jre1.8.0_241"
C:\Hadoop\Java <<===>> C:\Program Files\Java\jre1.8.0_241에 대한 교차점을 만들었습니다.
```

➢ **Finally, you can see JAVA folder emerges in your "C:\Hadoop"**

# Configure PATH

➢ **Enter Control Panel > System and Security > System > Advanced System Settings > Environment variables**

➢ **Then, Add variable "HADOOP_HOME" and "JAVA_HOME" like following**



*HADOOP_HOME :*
The folder where Hadoop files installed



*JAVA_HOME :*
The link folder to JAVA

# Configure PATH

➢ **Add "bin" folder of Hadoop folder to "Path" variable**

*Bin folder directory: YOUR_SPARK_PATH\Hadoop\bin*

➢ **For example, look at the following**

# Hadoop Configuration

➢ **Before start Hadoop file system, we must edit some configuration file**

- **Open %HADOOP_HOME%\etc\Hadoop\core-site.xml**



- **Edit CONFIGURATION like following**

```
<configuration>
 <property>
  <name>fs.defaultFS</name>
      <value>hdfs://localhost:9000</value>
 </property>
</configuration>
```



28

# Hadoop Configuration

➢ **Before start Hadoop file system, we must edit some configuration file**

- **Open %HADOOP_HOME%\etc\Hadoop\hdfs-site.xml**

| 이름 | 수정한 날짜 |
|---|---|
| spark-2.4.5-bin-hadoop2.7 › Hadoop › etc › hadoop | |
| hadoop-policy | 2015-08-05 오 |
| hdfs-site | 2020-03-31 오 |
| httpfs-env | 2015-08-05 오 |

```
See the License for the specific language governing p
limitations under the License. See accompanying LICEN
-->

<!-- Put site-specific property overrides in this file. -->

<configuration>

</configuration>
```

- **Edit CONFIGURATION like following**

```
<configuration>
 <property> <name>dfs.replication</name> <value>1</value> </property>
 <property>
  <name>dfs.namenode.name.dir</name> <value>file:/hadoop/data/dfs/namenode</value>
 </property>
 <property>
  <name>dfs.datanode.data.dir</name> <value>file:/hadoop/data/dfs/datanode</value>
 </property>
</configuration>
```

```
<configuration>
<property>
<name>dfs.replication</name>
<value>1</value>
</property>
<property>
<name>dfs.namenode.name.dir</name>
<value>file:/hadoop/data/dfs/namenode</value>
</property>
<property>
<name>dfs.datanode.data.dir</name>
<value>file:/hadoop/data/dfs/datanode</value>
</property>
</configuration>
```

29

# Hadoop Configuration

➢ **Before start Hadoop file system, we must edit some configuration file**

- **Open %HADOOP_HOME%\etc\Hadoop\yarn-site.xml**

e › spark-2.4.5-bin-hadoop2.7 › Hadoop › etc › hadoop

| 이름 ^ | 수정한 날짜 |
|---|---|
| yarn-env | 2015-08-05 오 |
| yarn-env | 2015-08-05 오 |
| yarn-site | 2020-03-31 오 |

```
<configuration>

<!-- Site specific YARN configuration properties -->

</configuration>
```

- **Edit CONFIGURATION like following**

```
<configuration>
<!-- Site specific YARN configuration properties -->
<property><name>yarn.nodemanager.aux-services</name><value>mapreduce_shuffle</value></property>
<property>
<name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
<value>org.apache.hadoop.mapred.ShuffleHandler</value>
</property>
<property><name>yarn.application.classpath</name>
<value>
%HADOOP_HOME%\etc\hadoop,
%HADOOP_HOME%\share\hadoop\common\*,
%HADOOP_HOME%\share\hadoop\common\lib\*,
%HADOOP_HOME%\share\hadoop\mapreduce\*,
%HADOOP_HOME%\share\hadoop\mapreduce\lib\*,
%HADOOP_HOME%\share\hadoop\hdfs\*,
%HADOOP_HOME%\share\hadoop\hdfs\lib\*,
%HADOOP_HOME%\share\hadoop\yarn\*,
%HADOOP_HOME%\share\hadoop\yarn\lib\*
</value></property></configuration>
```

```
<configuration>
<!-- Site specific YARN configuration properties -->
<property>
<name>yarn.nodemanager.aux-services</name>
<value>mapreduce_shuffle</value> </property>
<property>
<name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
<value>org.apache.hadoop.mapred.ShuffleHandler</value>
</property>
<property> <name>yarn.application.classpath</name>
<value>
%HADOOP_HOME%\etc\hadoop,
%HADOOP_HOME%\share\hadoop\common\*,
%HADOOP_HOME%\share\hadoop\common\lib\*,
%HADOOP_HOME%\share\hadoop\mapreduce\*,
%HADOOP_HOME%\share\hadoop\mapreduce\lib\*,
%HADOOP_HOME%\share\hadoop\hdfs\*,
%HADOOP_HOME%\share\hadoop\hdfs\lib\*,
%HADOOP_HOME%\share\hadoop\yarn\*,
%HADOOP_HOME%\share\hadoop\yarn\lib\*
</value>
</property>
</configuration>
```
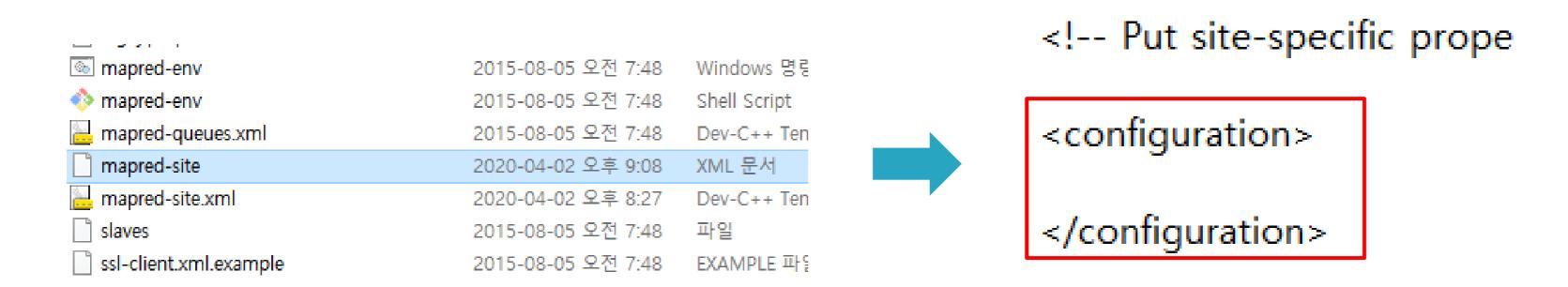
30

# Hadoop Configuration

➢ **Before start Hadoop file system, we must edit some configuration file**

- **Open %HADOOP_HOME%\etc\Hadoop\mapred-site.xml.template**

- **Make new file "mapred-site.xml" and copy the contents of mapred-site.xml.template**

| | | | |
|---|---|---|---|
| mapred-env | 2015-08-05 오전 7:48 | Windows 명령 | |
| mapred-env | 2015-08-05 오전 7:48 | Shell Script | |
| mapred-queues.xml | 2015-08-05 오전 7:48 | Dev-C++ Ten | |
| mapred-site | 2020-04-02 오후 9:08 | XML 문서 | |
| mapred-site.xml | 2020-04-02 오후 8:27 | Dev-C++ Ten | |
| slaves | 2015-08-05 오전 7:48 | 파일 | |
| ssl-client.xml.example | 2015-08-05 오전 7:48 | EXAMPLE 파일 | |

```
<!-- Put site-specific prope

<configuration>

</configuration>
```

- **Edit CONFIGURATION in mapred-site.xml file like following**

```
<configuration>
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>
</configuration>
```

```
<configuration>
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>
</configuration>
```

# Hadoop Configuration

- At last, we can use Hadoop file system

- If you want to Format Namenode, then input like following

   *%HADOOP_HOME%\bin\hdfs namenode -format*

- For example,

```
(base) C:\Users\BigDataLab>%HADOOP_HOME%\bin\hdfs namenode -format
20/03/30 17:11:29 INFO namenode.NameNode: STARTUP_MSG:
/************************************************************
STARTUP_MSG: Starting NameNode
STARTUP_MSG:    host = DESKTOP-4FSOHKP/192.168.56.1
STARTUP_MSG:    args = [-format]
STARTUP_MSG:    version = 2.7.1
STARTUP_MSG:    classpath = C:\Users\BigDataLab\spark-2.4.5-bin-hadoop2.7\Hadoop\etc\hadoop;
```
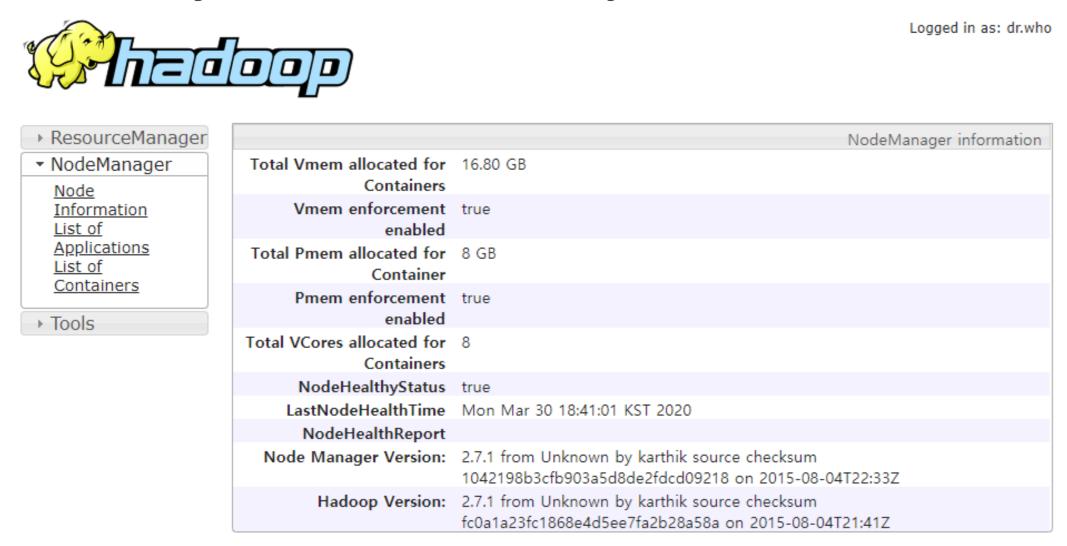
# Start Hadoop File System

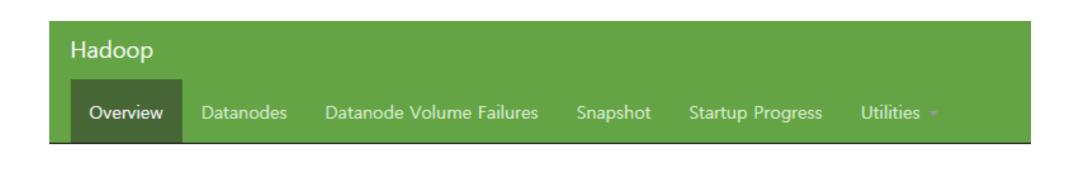➢ **Start DFS and YARN with the following command**

*%HADOOP_HOME%\sbin\start-dfs.cmd*

*%HADOOP_HOME%\sbin\start-yarn.cmd*

➢ **A few seconds later, enter the following links**

http://localhost:8088 & http://localhost:50070

➢ **Then you can see Hadoop websites.**

# Using Linux

# Virtual box installation

# Virtual Box

➢ **VirtualBox is a powerful virtualization product.**

# CentOS 7

➢ **Download CentOS 7 with the following link:**

http://mirror.kakao.com/centos/7.7.1908/isos/x86_64/

➢ **Download CentOS-x86_64-DVD-1908.iso**

# CentOS 7

➢ **Choose operation system in detail**

# CentOS 7

- ➢ **Minimum volume of the hard disk**

  - You have to make volume of the hard disk more than 15 GB.
  - Unless, you are able to get out of disk.

# CentOS 7

➢ **CentOS installation process**

# CentOS 7

➢ **CentOS installation process**

# CentOS 7

➢ **CentOS installation process**

# CentOS 7

➢ **CentOS installation process**

# CentOS 7

➢ **CentOS installation process**

# CentOS 7

➢ **CentOS installation process**

# CentOS 7

➢ **CentOS installation process**

# CentOS 7

➢ **CentOS installation process**

# Download Python3

➢ **Start CentOS, then move into the terminal**

➢ **Install python3 with the following command**

- Add repository to yum

    sudo yum install –y https://centos7.iuscommunity.org/ius-release.rpm

- Install the library

    sudo yum install –y python36u python36-libs python36-devel python36u-pip

- Modify Alias

    sudo unlink /bin/python

    sudo ln -s /bin/python3.6  /bin/python

# Download Pip

➢ **Install Pip which is a tool for downloading python library**

➢ **Install Pip with the following command**

- sudo curl https://bootstrap.pypa.io/get-pip.py -o get-pip.py
- python get-pip.py

- **You can check the pip version using "pip --version" command**

```
[wj-lee@localhost ~]$ pip --version
pip 20.0.2 from /usr/local/lib/python3.6/site-packages/pip (python 3.6)
```

# Download Packages

➢ **Install libraries using pip3 with the following command:**

  pip3 install numpy
  pip3 install scikit-learn
  pip3 install matplotlib

➢ **Install pyspark library with the following command:**

  pip3 install pyspark –U --no-cache

# Spark installation

# Download Spark

> **Go http://spark.apache.org/ and click "Download"**

# Download Spark

➢ **Choose spark 2.4.5 & Apache Hadoop 2.7**



➢ **Then, download "spark-2.4.5-bin-hadoop2.7.tgz"**

➢ **Or, you can download the file here:**

https://www.apache.org/dyn/closer.lua/spark/spark-2.4.5/spark-2.4.5-bin-hadoop2.7.tgz

http://spark.apache.org/

# Download Spark

➢ **Unpack your .tgz file**

- Move *spark-2.4.5-bin-hadoop2.7.tgz* file to your $HOME directory
  **sudo mv spark-2.4.5-bin-hadoop2.7 $HOME (at your current directory where this file installed)**

- Unpack your .tgz file like following command:
  **tar zxvf spark-2.4.5-bin-hadoop2.7.tgz**

➢ **Set the path**
  **export SPARK_HOME=$HOME/spark-2.4.5-bin-hadoop2.7**
  **export PATH=$PATH:$SPARK_HOME/bin**
  **echo 'export SPARK_HOME=$HOME/spark-2.4.5-bin-hadoop2.7' >> .bash_profile**
  **echo 'export PATH=$PATH:$SPARK_HOME/bin' >> .bash_profile**
  **echo 'export SPARK_HOME= $HOME/spark-2.4.5-bin-hadoop2.7' >> ~/.bashrc**
  **echo 'export PATH=$PATH:$SPARK_HOME/bin' >> ~/.bashrc**

https://gist.github.com/darcyliu/d47edccb923b0f03280a4cf8b66227c1

# Start Pyspark in CentOS

> ➢ **You can execute Pyspark entering just "pyspark" at command line.**

> ➢ **Then, you can import pyspark here.**

```
[root@localhost practice1]# pyspark
Python 3.6.8 (default, Aug  7 2019, 17:28:10)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-39)] on linux
Type "help", "copyright", "credits" or "license" for more information.
20/03/26 01:22:51 WARN Utils: Your hostname, localhost.localdomain resolves to a loopback address: 127.0.0.1; using 10.0.2.15
interface enp0s3)
20/03/26 01:22:51 WARN Utils: Set SPARK_LOCAL_IP if you need to bind to another address
20/03/26 01:22:52 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes
able
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
20/03/26 01:22:57 WARN Utils: Service 'SparkUI' could not bind on port 4040. Attempting port 4041.
20/03/26 01:22:57 WARN Utils: Service 'SparkUI' could not bind on port 4041. Attempting port 4042.
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 2.4.5
      /_/

Using Python version 3.6.8 (default, Aug  7 2019 17:28:10)
SparkSession available as 'spark'.
>>> import pyspark
>>>
```

# Hadoop installation

# Install Hadoop

- ➢ **Check whether JAVA is already installed : java -version**
- ➢ **Install ssh and and pdsh**
    - **sudo yum install openssh-server openssh-clients openssh-askpass**
    - **sudo yum install pdsh**
- ➢ **In case of yum exception, edit the first line of configuration files like below**
    - **sudo vi /usr/bin/yum**
    - **sudo vi /usr/libexec/urlgrabber-ext-down**

### #! /usr/bin/python ➡ #! /usr/bin/python2.7

- ❖ *In Vim editior, click "i" for editing(edit mode)*
- ❖ *click ESC(command mode) and ":wq" for save and exit.*
- ❖ *If you want to exit without save, click ":q!" and type "q!".*

# Install Hadoop

➤ **Go to Hadoop website and download package**

**https://hadoop.apache.org/release/2.7.7.html**

➤ **Click "Download tar.gz"**

# Install Hadoop

- ➢ **Move tar.gz package from Downloads directory to home directory**
- ➢ **Unpack tar.gz package and check new directory for Hadoop**

<div align="center">

**tar zxvf hadoop-2.7.7.tar.gz**

**ls hadoop-2.7.7**

</div>

```
[bigdatalab@localhost ~]$ ls hadoop-2.7.7
LICENSE.txt   NOTICE.txt   README.txt   bin   etc   include   lib   libexec   sbin   share
```

- ➢ **Path configuration**

    **export HADOOP_HOME=$HOME/hadoop-2.7.7**

    **export PATH=$PATH:$HADOOP_HOME/bin**

    **echo 'export HADOOP_HOME=$HOME/hadoop-2.7.7' >> .bash_profile**

    **echo 'export PATH=$PATH:$HADOOP_HOME/bin' >> .bash_profile**

    **echo 'export HADOOP_HOME=$HOME/hadoop-2.7.7' >> ~/.bashrc**

    **echo 'export PATH=$PATH:$HADOOP_HOME/bin' >> ~/.bashrc**

# Hadoop Configuration

➢ **Edit JAVA path configuration file**

   **sudo vi $HADOOP_HOME/etc/hadoop/hadoop-env.sh**

   **edit "export JAVA_HOME=${JAVA_HOME}" to "export JAVA_HOME=/usr/lib/jvm/jre-1.8.0-openjdk"**

```
# Licensed to the Apache Software Foundation (ASF) under one
# or more contributor license agreements.  See the NOTICE file
# distributed with this work for additional information
# regarding copyright ownership.  The ASF licenses this file
# to you under the Apache License, Version 2.0 (the
# "License"); you may not use this file except in compliance
# with the License.  You may obtain a copy of the License at
#
#     http://www.apache.org/licenses/LICENSE-2.0
#
# Unless required by applicable law or agreed to in writing, software
# distributed under the License is distributed on an "AS IS" BASIS,
# WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
# See the License for the specific language governing permissions and
# limitations under the License.


# Set Hadoop-specific environment variables here.

# The only required environment variable is JAVA_HOME.  All others are
# optional.  When running a distributed configuration it is best to
# set JAVA_HOME in this file, so that it is correctly defined on
# remote nodes.

# The java implementation to use.
export JAVA_HOME=/usr/lib/jvm/jre-1.8.0-openjdk
```

❖ *In Vim editior, click "i" for editing(edit mode)*
❖ *click ESC(command mode) and ":wq" for save and exit.*
❖ *If you want to exit without save, click ":q!" and type "q!".*

# Hadoop Configuration

➢ **Before start Hadoop file system, we must edit some configuration files**

- **sudo vi $HADOOP_HOME/etc/hadoop/core-stie.xml**
- **Edit CONFIGURATION like following**

```
<configuration>
  <property>
    <name>fs.defaultFS</name>
    <value>hdfs://localhost:9000</value>
  </property>
</configuration>
```

```
<!-- Put site-specific property overrides in this file. -->

<configuration>
        <property>
                <name>fs.defaultFS</name>
                <value>hdfs://localhost:9000</value>
        </property>
</configuration>
~
~
~
~
~
```
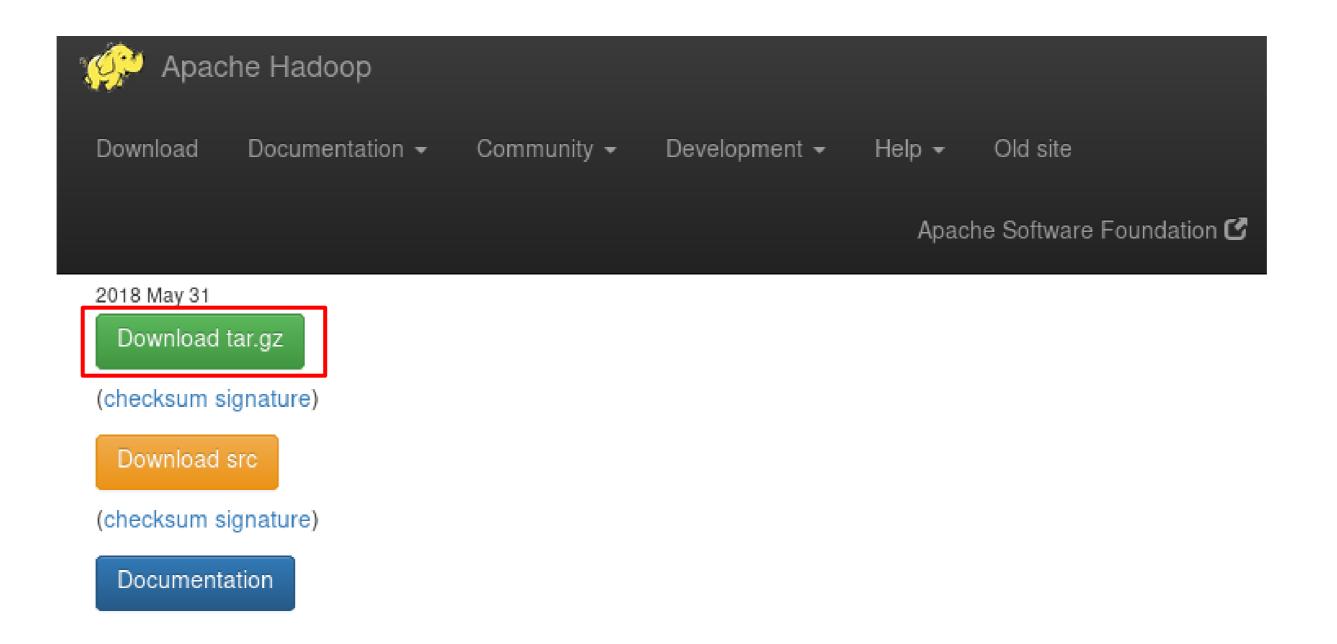
❖ *In Vim editior, click "i" for editing(edit mode)*
❖ *click ESC(command mode) and ":wq" for save and exit.*
❖ *If you want to exit without save, click ":q!" and type "q!".*

# Hadoop Configuration

➤ **Before start Hadoop file system, we must edit some configuration files**

  • **sudo vi $HADOOP_HOME/etc/hadoop/hdfs-site.xml**

  • **Edit CONFIGURATION like following**

```
<configuration>
    <property>
        <name>dfs.replication</name>
<value>1</value>
</property>
<property>
        <name>dfs.namenode.name.dir</name>
        <value>file:/hadoop/data/dfs/namenode</value>
</property>
<property>
        <name>dfs.datanode.data.dir</name>
        <value>file:/hadoop/data/dfs/datanode</value>
</property>
</configuration>
```

```
<configuration>
        <property>
                <name>dfs.replication</name>
                <value>1</value>
        </property>
        <property>
                <name>dfs.namenode.name.dir</name>
                <value>file:/hadoop/data/dfs/namenode</value>
        </property>
        <property>
                <name>dfs.datanode.data.dir</name>
                <value>file:/hadoop/data/dfs/datanode</value>
        </property>
</configuration>
~
```

❖ *In Vim editior, click "i" for editing(edit mode)*
❖ *click ESC(command mode) and ":wq" for save and exit.*
❖ *If you want to exit without save, click ":q!" and type "q!".*

# Hadoop Configuration

➢ **Before start Hadoop file system, we must edit some configuration files**

• **sudo vi $HADOOP_HOME/etc/hadoop/mapred-site.xml**

• **Edit CONFIGURATION like following**

*<configuration>*
  *<property>*
    *<name>mapreduce.framework.name</name>*
    *<value>yarn</value>*
  *</property>*
  *<property>*
    *<name>mapreduce.application.classpath</name>*

*<value>$HADOOP_MAPRED_HOME/share/hadoop/mapreduce/*:$HADOOP_MAPRED_HOME/share/hadoop/mapreduce/lib/*</value>*
  *</property>*
*</configuration>*

```
<configuration>
        <property>
                <name>mapreduce.framework.name</name>
                <value>yarn</value>
        </property>
        <property>
                <name>mapreduce.application.classpath</name>
                <value>$HADOOP_MAPRED_HOME/share/hadoop/mapreduce/*:$HADOOP_MAPRED_
        </property>
</configuration>
~
~
~
~
~
~
```

❖ *In Vim editior, click "i" for editing(edit mode)*
❖ *click ESC(command mode) and ":wq" for save and exit.*
❖ *If you want to exit without save, click ":q!" and type "q!".*

# Hadoop Configuration

➤ **Before start Hadoop file system, we must edit some configuration files**

- sudo vi $HADOOP_HOME/etc/hadoop/yarn-site.xml
- Edit CONFIGURATION like following

*<configuration>*
  *<property>*
    *<name>yarn.nodemanager.aux-services</name>*
    *<value>mapreduce_shuffle</value>*
  *</property>*
  *<property>*
    *<name>yarn.nodemanager.env-whitelist</name>*

*<value>JAVA_HOME,HADOOP_COMMON_HOME,HADO OP_HDFS_HOME,HADOOP_CONF_DIR,CLASSPATH_PRE PEND_DISTCACHE,HADOOP_YARN_HOME,HADOOP_M APRED_HOME</value>*
    *</property>*
*</configuration>*

```
<configuration>
        <property>
                <name>yarn.nodemanager.aux-services</name>
                <value>mapreduce_shuffle</value>
        </property>
        <property>
                <name>yarn.nodemanager.env-whitelist</name>
                <value> JAVA_HOME,HADOOP_COMMON_HOME,HADOOP_HDFS_HOME,HADOOP_CONF_DIR,C
LASSPATH_PREPEND_DISTCHACHE,HADOOP_YARN_HOME,HADOOP_MAPRED_HOME</value>
        </property>
</configuration>
```

❖ *In Vim editior, click "i" for editing(edit mode)*
❖ *click ESC(command mode) and ":wq" for save and exit.*
❖ *If you want to exit without save, click ":q!" and type "q!".*

64

# Hadoop Configuration

➢ **To format Namenode enter the command: hdfs namenode -format**

```
20/03/31 18:41:24 INFO util.ExitUtil: Exiting with status 0
20/03/31 18:41:24 INFO namenode.NameNode: SHUTDOWN_MSG:
/************************************************************
SHUTDOWN_MSG: Shutting down NameNode at localhost/127.0.0.1
************************************************************/
[bigdatalab@localhost ~]$
```

➢ **Start DFS and YARN (You need to type password multiple times)**

**sudo $HADOOP_HOME/sbin/start-dfs.sh**

**sudo $HADOOP_HOME/sbin/start-yarn.sh**

➢ **Check "http://localhost:8088" and "http://localhost:50070"**

# Spark Implementation

# Steps for Spark Implementation

```
from pyspark import SparkConf, SparkContext
conf = SparkConf()
conf.set("spark.master", "local")
conf.set("spark.app.name", "My app")
sc = SparkContext(conf=conf)
```

```
>>> lines = sc.parallelize(["pandas", "i like pandas"])
>>> print(lines.first())
pandas
...
>>> lines = sc.textFile("test.txt")
>>> print(lines.first())
I like the Spark
```

1. **Initialize a SparkContext**
   - ➤ Import the Spark Package in your program
   - ➤ Configure Spark with SparkConf
   - ➤ Call *set* to add configuration values

2. **Create RDDs (Import data)**
   - ➤ parallelize()
   - ➤ textFile()
     - ▪ load data from an external storage

# Steps for Spark Implementation

## 3. RDD Operations

RDDs support two types of operations, *transformation* and *action*

1) Transformation
   - Operations on RDDs that return a new RDD
   - Transformed RDDs are computed only when you use them with an action function

2) Action
   - Operations that return a final value to the driver program or write data to an external storage system
   - Actions force the evaluation of the transformations required for the RDD they were called on

* ref.) https://spark.apache.org/docs/latest/rdd-programming-guide.htm

❖ Passing Function
   - Most of Spark's transformations, and some of its actions, depend on passing in function that are used to compute data
   - In python, we can pass in lambda expressions, top-level functions, or locally defined

```
>>> lines = sc.parallelize(["error in python", "pandas", "error in spark"])
>>> errorRDD = lines.filter(lambda x : "error" in x)
>>> for line in errorRDD.collect():
...     print(line)
...
error in python
error in spark
```

```
>>> input = sc.textFile("test.txt")
>>> print(input.count())
4
>>> for line in input.take(2):
...     print(line)
...
I like the Spark
It's rainy day
```

```
upper_input = input.map(lambda x: x.upper())
```

```
>>> def UPPER(x):
...     return x.upper()
...
>>> upper_input = input.map(UPPER)
```

68

# Steps for Spark Implementation

**4. Use the Spark in an application**

&#10148; In python, you simply write applications as Python scripts

&#10148; You must run them using the

   *$ spark-submit ScriptName.py*

&#10148; If you run your program on single machine with specific cores, use

   *$ spark-submit --master local[N] SctriptName.py*

   \* N : the number of cores you want to execute

# Issues

# Hadoop Issues

➢ **ERROR : Output directory already exists**

```
20/04/01 16:23:07 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/04/01 16:23:07 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/04/01 16:23:07 ERROR streaming.StreamJob: Error Launching job : Output directory hdfs://localhost:9000/output alre
ady exists
Streaming Command Failed!
```

1.  **Delete HDFS output directory using below command**

   **(Windows) hdfs dfs -rm /YOUR_DIRECTORY/***
   **hdfs dfs -rmdir /YOUR_DIRECTORY/**
   **(Linux) sudo %HADOOP_HOME%/bin/hdfs dfs –rm /YOUR_DIRECTORY/***
   **sudo %HADOOP_HOME%/bin/hdfs dfs -rmdir /YOUR_DIRECTORY/**

2.  **Try again**

# Hadoop Issues

➢ **ERROR : Initialization failed for Block pool**

```
20/04/01 16:42:46 FATAL datanode.DataNode: Initialization failed for Block pool <registering> (Datanode Uuid unassigned)
 service to localhost/127.0.0.1:9000. Exiting.
java.io.IOException: All specified directories are failed to load.
```

1. **Delete "data" folder.**

    **(Windows) rm -r C:\Hadoop\data OR delete directly**

    **(Linux) sudo rm -r  /hadoop/data**


2. **Then, format namenode**

    **(Windows)hdfs namenode -format**

    **(Linux) sudo $HADOOP_HOME/bin/hdfs namenode -format**
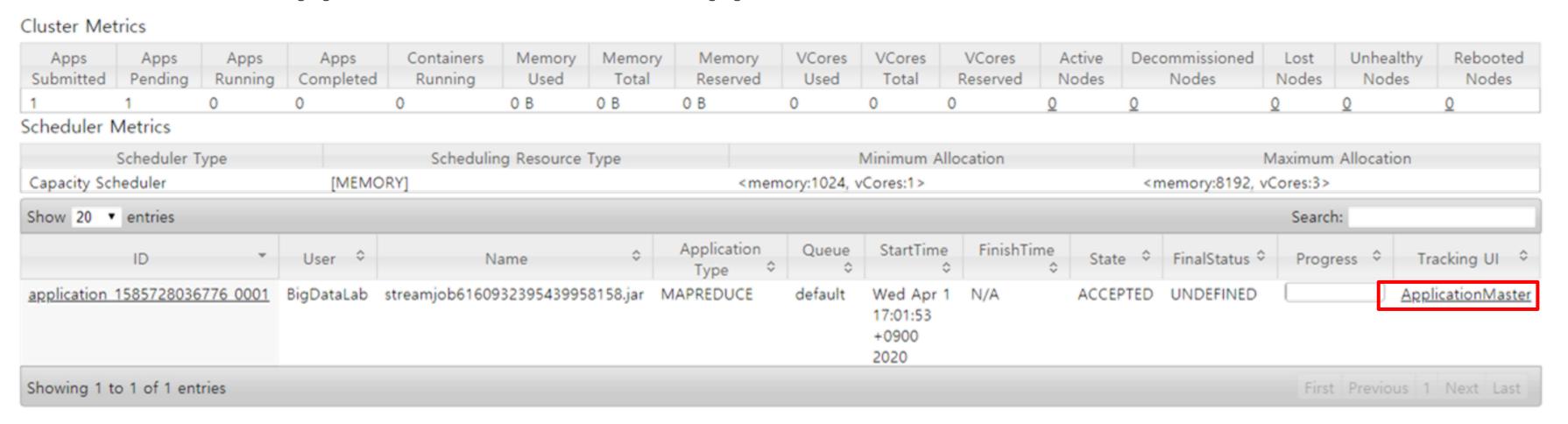

3. **Try again.**

# Hadoop Issues

➤ **ERROR : Process stopped more than a few minutes before starting Mapreduce**

```
20/04/01 17:01:50 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/04/01 17:01:50 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
20/04/01 17:01:52 INFO mapred.FileInputFormat: Total input paths to process : 1
20/04/01 17:01:52 INFO mapreduce.JobSubmitter: number of splits:3
20/04/01 17:01:53 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1585728036776_0001
20/04/01 17:01:53 INFO impl.YarnClientImpl: Submitted application application_1585728036776_0001
20/04/01 17:01:53 INFO mapreduce.Job: The url to track the job: http://DESKTOP-4FS0HKP:8088/proxy/application_1585728
036776_0001/
20/04/01 17:01:53 INFO mapreduce.Job: Running job: job_1585728036776_0001
```
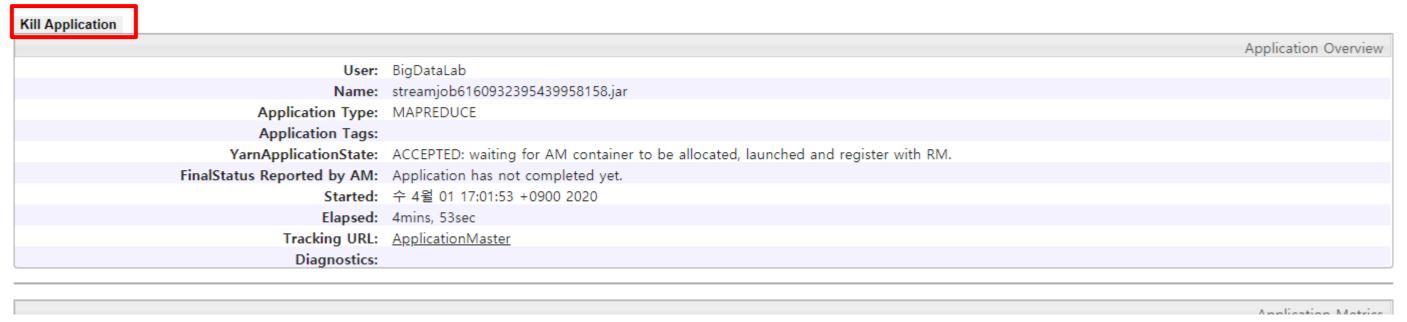
1. **Go to Hadoop cluster admin page**

   **(Windows, Linux) http://localhost:8088/cluster**

# Hadoop Issues

2. **Find current application and click ApplicationMaster**



3. **Kill application and try again**



4. **If same issue repeated, format namenode and try again**

# Hadoop Issues

➢ **ERROR : SAFE MODE**

   **When dfs and yarn closed without command (sbin/stop-all),**

   **namenode and datanode can be safe mode which inserting or deleting file is not allowed.**

1. **Break safe mode using command**

   **(Windows) hadoop dfsadmin -safemode leave**

   **(Linux) sudo $HADOOP_HOME/bin/hadoop dfsadmin -safemode leave**

※ **Use stop-all command when closing dfs and yarn**

   **(Windows) %HADOOP_HOME%\sbin\stop-all.cmd**

   **(Linux) sudo $HADOOP_HOME/sbin/stop-all.sh**