# Installation of Apache Hadoop on Windows

**Step 1**: Download Hadoop, Java and winutils from below links.

- We're going to install Hadoop version 3.1.0
  - [Hadoop 3.1.0](#)
  - [Java](#)
  - [Winutils](#)
- Extract all three files (Hadoop will take a while).
- Run the Java installer but change the destination folder from the default "C:\Program Files\Java\jdk.1.8.0_191" to just "C:\Java".
- Make a directory C:\Hadoop
- Move all contents of Hadoop-3.1.0 to C:\Hadoop

**Step 2**: Setup Environment variables

- Go to Control Panel > System and Security > System > Advanced System Settings > Environment Variables
- Add new System variables (bottom box) called:
  - JAVA_HOME --> C:\Java
  - HADOOP_HOME --> C:\Hadoop
- Edit **Path** and add the following:
  - C:\ Java\jdk1.8.0_111\bin
  - C:\Hadoop\bin
- Check in by opening CMD ( Win + R -> cmd ):
  - Java -version
  - Hdfs – version

**Step 3**: Setup Hadoop Configurations

- Go to C:\Hadoop\etc\hadoop and edit\create core-site.xml

```
1. <configuration>
2. <property>
3. <name>fs.defaultFS</name>
4. <value>hdfs://localhost:9000</value>
5. </property>
6. </configuration>
```

- In the same directory, edit (or create) mapred-site.xml with the following contents:

```
1. <configuration>
2.   <property>
3.     <name>mapreduce.framework.name</name>
4.     <value>yarn</value>
```

```
5.    </property>
6. </configuration>
```

- Next, edit (or create) hdfs-site.xml:

```
1.  <configuration>
2.    <property>
3.      <name>dfs.replication</name>
4.      <value>1</value>
5.    </property>
6.    <property>
7.      <name>dfs.namenode.name.dir</name>
8.      <value>file:///C:/hadoop/namenode</value>
9.    </property>
10.   <property>
11.       <name>dfs.datanode.data.dir</name>
12.       <value>file:///C:/hadoop/datanode</value>
13.   </property>
14. </configuration>
```

- Finally, edit yarn-site.xml

```
1.  <configuration>
2.    <property>
3.      <name>yarn.nodemanager.aux-services</name>
4.      <value>mapreduce_shuffle</value>
5.    </property>
6.    <property>
7.      <name>yarn.nodemanager.auxservices.mapreduce.shuffle.class</name>
8.      <value>org.apache.hadoop.mapred.ShuffleHandler</value>
9.    </property>
10. </configuration>
```

- The last thing we need to do is create the directories that we referenced in hdfs-site.xml (DIY).

**Step 4**: Patch Hadoop

- Hadoop was mainly created for Unix, so we need to patch it to make it work on windows.

- Just copy bin folder from winutils and paste it in C:\hadoop.
- If it asks for overwrite then allow it.
- copy **hadoop-yarn-server-timelineservice-3.0.3** from C:\hadoop\share\hadoop\yarn\timelineservice to C: \hadoop\share\hadoop\yarn (the parent directory).

Step 5: Run HDFS

- Open CMD and run **hdfs namenode -format**
- Finally, you can boot HDFS by running start-dfs.cmd and start-yarn.cmd in cmd. (It will open 4 cmd windows).
- You can monitor these windows by typing "**jps"** in cmd.
- You can see list of all applications by opening http://localhost:8088

Congratulations, if you are here without any error (or by solving many errors) you've successfully installed Hadoop cluster in your windows machine.

Some common problems:

- Muhammad Bilal Yar's Hadoop 2.8.0 walkthrough

- java.net.URISyntaxException

- java.lang.UnsatisfiedLinkError

- FATAL resourcemanager.ResourceManager

- localhost:50070 error

- Kuldeep Singh's walkthrough and troubleshooting guide

- Jacek Laskowski's GitBook

- java.io.IOException: Incompatible clusterIDs

- HDFS basic commands

- Spark basic commands

Reference : https://dev.to/awwsmm/installing-and-running-hadoop-and-spark-on-windows-33kc

**Running your Program with python:**

- Start HDFS server : `start-all.cmd`
- Format Namenode : `hdfs namenode -format`
- Copy data file to HDFS : `hdfs dfs -put data.txt /`
- List files (same as unix) : `hdfs dfs -ls /`
- Run MapReduce Job :
    - `hadoop jar C:\hadoop\share\hadoop\tools\lib\hadoop-streaming-3.1.0.jar -file path\to\mapper.py -mapper "python mapper.py" -file path\to\reducer.py -reducer "python reducer.py" -input text.txt -output out`
- It will save file in HDFS and you must bring it back via:

    - `hdfs dfs -get /user/hadoop/file localfile`