

## Homework #2: Exploring HDFS Metadata Using XML & XPath

Due: February 18, Friday (11:59pm)

100 points

In this homework, we will explore the metadata stored in the namenode of HDFS. You can obtain such metadata by using the Offline Image Viewer (oiv) tool provided by Hadoop (<https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-hdfs/HdfsImageViewer.html>).

For example,

```
<Your Hadoop-installation-dir>/bin/hdfs oiv -i /tmp/hadoop-ec2-user/dfs/name/current/
fsimage_00000000000000000564 -o fsimage564.xml -p XML
```

will export the metadata stored in the specified fsimage (file system image) to an XML file called fsimage546.xml.

```
▼ <INodeSection>
  <lastInodeId>16422</lastInodeId>
  <numInodes>38</numInodes>
  ▼ <inode>
    <id>16385</id>
    <type>DIRECTORY</type>
    <name/>
    <mtime>1581231015982</mtime>
    <permission>ec2-user:supergroup:0755</permission>
    <nsquota>9223372036854775807</nsquota>
    <dsquota>-1</dsquota>
  </inode>
  ▼ <inode>
    <id>16386</id>
    <type>DIRECTORY</type>
    <name>user</name>
    <mtime>1581231034866</mtime>
    <permission>ec2-user:supergroup:0755</permission>
    <nsquota>-1</nsquota>
    <dsquota>-1</dsquota>
  </inode>
  ▼ <inode>
    <id>16387</id>
    <type>DIRECTORY</type>
    <name>ec2-user</name>
    <mtime>1581875598912</mtime>
    <permission>ec2-user:supergroup:0755</permission>
    <nsquota>-1</nsquota>
    <dsquota>-1</dsquota>
  </inode>
  ▼ <INodeDirectorySection>
    ▼ <directory>
      <parent>16385</parent>
      <child>16386</child>
    </directory>
    ▼ <directory>
      <parent>16386</parent>
      <child>16387</child>
    </directory>
    ▼ <directory>
      <parent>16387</parent>
      <child>16390</child>
      <child>16412</child>
      <child>16401</child>
      <child>16391</child>
      <child>16388</child>
    </directory>
    ▼ <directory>
      <parent>16388</parent>
      <child>16389</child>
    </directory>
    ▼ <directory>
      <parent>16391</parent>
      <child>16392</child>
      <child>16393</child>
      <child>16394</child>
      <child>16395</child>
      <child>16396</child>
      <child>16397</child>
      <child>16398</child>
      <child>16399</child>
      <child>16400</child>
    </directory>
```

Fsimage has a INodeSection listing metadata about each inode and a INodeDirectorySection describing the directory structure, as show above. Note that id of inode is its inumber; and the directory nodes are represented by their inumbers, e.g., 16385.

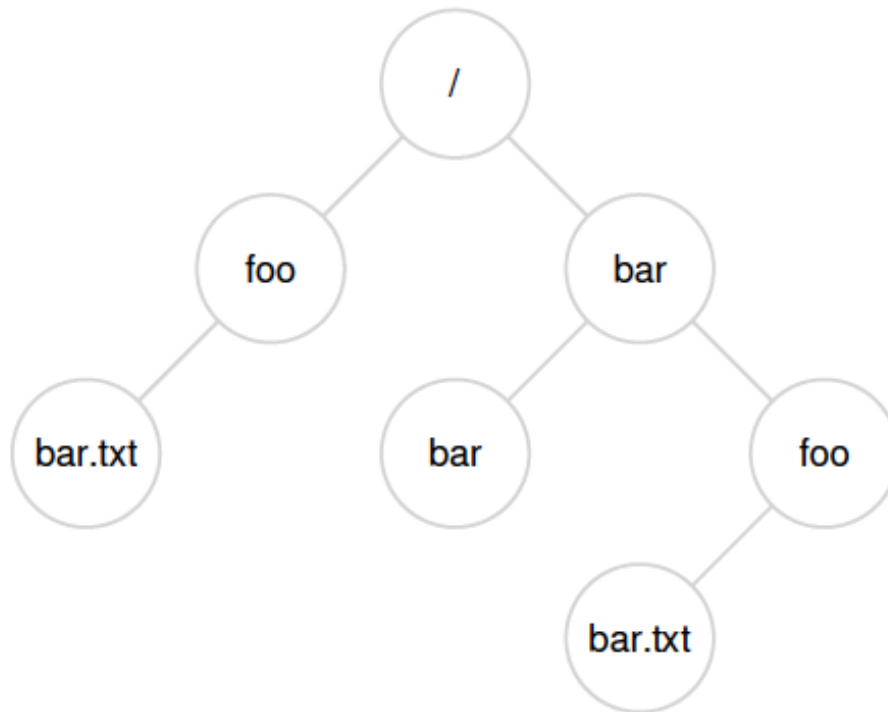
## DSCI 551 – Spring 2022

Your task is to implement a Python program stats.py that takes a fsimage file in XML and outputs an JSON file that contains the following statistics about the file system.

{“number of files”: 5, “number of directories”: 10, “maximum depth of directory tree”: 4

“file size”: {“max”: 3518, “min”: 16}}

Note the maximum depth of directory tree is the number of levels of the tree, e.g., the maximum depth of the following directory tree is 4. Note if the file system does not contain any files, then you should not output the statistics about the “file size”.



Permitted libraries: lxml.

Execution format: python3 stats.py <fsimage xml file> <stats json file>

e.g.,

python3 stats.py fsimage564.xml stats.json

**Submission:** submit stats.py

## DSCI 551 – Spring 2022

### Checklist for Submission :

1. **DO NOT return anything we didn't ask for.** For example, "please enter XXX: \_\_\_\_". Please no. We have given you the EXACT output format. Please just follow them.
2. Use **ONLY relative path**. You should assume your scripts will be run in the directory where the scripts are at. For example, no 'C:\\homework1\\...' or '/Users/ blabla/....'
3. Make sure that you are able to run the code according to the execution format mentioned above in the questions.
4. Double-check your files before submitting them. Please use python3 to complete the homework and try to maintain the python version as 3.7. Do not use any libraries other than the ones specified in the handout. You can use EC2 to test your code, and python 3.7 is preinstalled on EC2.
5. You can submit it multiple times on DEN but only the latest attempt will be graded.
6. Please only submit the stats.py file, and do not include your output file in the submission.