800+ Q&As    |    Logout    |    Contact

# Java-Success.com

Prepare to fast-track, choose & go places with 800+ Java & Big Data Q&As with lots of code & diagrams.

search here …        Go

Home    Why? ▾    300+ Java FAQs ▾    300+ Big Data FAQs ▾    Courses ▾

👤 Membership ▾    Your Career ▾

Home › bigdata-success.com › Tutorials - Big Data › TUT - File Formats › 01a: Convert XML file To Sequence File – writing & reading – Local File System

# 01a: Convert XML file To Sequence File – writing & reading – Local File System

📅 Posted on May 1, 2016

Sequence files are good for saving **raw data** into HDFS. Sequence files are **compressible** and **splittable**. It is also useful for combining a number of smaller files into a single say 64MB or larger sequence file as HDFS is more suited for larger files. Text files like csv, xml, json, etc can be stored on HDFS in **sequence file format**.

Step 1: Create a Java project using Maven.

## 300+ Java Interview FAQs

300+ Java FAQs 🔥                      ⌄

16+ Java Key Areas Q&As              ⌄

150+ Java Architect FAQs             ⌄

80+ Java Code Quality Q&As           ⌄

150+ Java Coding Q&As                ⌄

## 300+ Big Data Interview FAQs

300+ Big Data FAQs 🔥                 ⌄

Tutorials - Big Data                 ›

   TUT - 🔷 Starting Big Data

   TUT - Starting Spark & Scala

```
1
2   mvn archetype:generate -DgroupId=com.mytutorial -D
3
```

**Step 2:** Import it into eclipse as a an "**existing maven project**".

**Step 3:** Add the hadoop dependencies to the **pom.xml** file.

```
1
2   <project xmlns="http://maven.apache.org/POM/4.0.0"
3       xsi:schemaLocation="http://maven.apache.org/P(
4       <modelVersion>4.0.0</modelVersion>
5       <groupId>com.mytutorial</groupId>
6       <artifactId>sequence-file</artifactId>
7       <packaging>jar</packaging>
8       <version>1.0-SNAPSHOT</version>
9       <name>sequence-file</name>
10      <url>http://maven.apache.org</url>
11
12      <properties>
13          <maven.compiler.source>1.8</maven.compiler
14          <maven.compiler.target>1.8</maven.compiler
15          <project.build.sourceEncoding>UTF-8</proje
16          <junit.version>4.8.1</junit.version>
17          <hadoop.version>2.7.2</hadoop.version>
18      </properties>
19
20      <dependencies>
21          <!-- JUnit -->
22          <dependency>
23              <groupId>junit</groupId>
24              <artifactId>junit</artifactId>
25              <version>${junit.version}</version>
26              <scope>test</scope>
27          </dependency>
28
29          <!-- Hadoop -->
30          <dependency>
31              <groupId>org.apache.hadoop</groupId>
32              <artifactId>hadoop-hdfs</artifactId>
33              <version>${hadoop.version}</version>
34              <exclusions>
35                  <exclusion>
36                      <groupId>javax.servlet</groupI
37                      <artifactId>*</artifactId>
38                  </exclusion>
39              </exclusions>
```

## 800+ Java Interview Q&As

300+ Core Java Q&As ⌄

300+ Enterprise Java Q&As ⌄

150+ Java Frameworks Q&As ⌄

120+ Companion Tech Q&As ⌄

Tutorials - Enterprise Java ⌄

```
40            </dependency>
41            <dependency>
42                <groupId>org.apache.hadoop</groupId>
43                <artifactId>hadoop-client</artifactId>
44                <version>${hadoop.version}</version>
45                <exclusions>
46                    <exclusion>
47                        <groupId>javax.servlet</groupl
48                        <artifactId>*</artifactId>
49                    </exclusion>
50                </exclusions>
51            </dependency>
52
53        </dependencies>
54
55 </project>
56
```

**Step 4**: Create **src/main/resources** source folder, and under that the "data" sub-folder. Create **"report.xml"**, which needs to be converted to a **sequence file format**.

```
1
2  <?xml version="1.0" encoding="UTF-8"?>
3  <transactionReports xmlns="http://mytutorial.com/
4      <transactionReport>
5          <report>
6              <reportNumber>9999</reportNumber>
7              <createdDatetime>2015-06-15T11:29:52+
8              <processedDatetime>2015-06-15T11:29:5
9              <reportStatusCode>Active</reportStatu:
10         </report>
11     </transactionReport>
12 </transactionReports>
13
```

**Step 5**: Write a stand-alone Java class to convert above report.xml file to report.seq. Write and then read. This example uses the **"LOCAL FILE SYSTEM"**

```
1
2  package com.mytutorial;
3
4  import java.io.File;
5  import java.io.IOException;
```

```java
 6  import java.net.URL;
 7  import java.util.List;
 8  import java.util.stream.Collectors;
 9
10  import org.apache.commons.io.FileUtils;
11  import org.apache.hadoop.conf.Configuration;
12  import org.apache.hadoop.fs.Path;
13  import org.apache.hadoop.io.BytesWritable;
14  import org.apache.hadoop.io.IOUtils;
15  import org.apache.hadoop.io.IntWritable;
16  import org.apache.hadoop.io.SequenceFile;
17  import org.apache.hadoop.util.ReflectionUtils;
18
19  public class ConvertXmlToSequence {
20      private static final String FILE_IN_PATH = "d
21      private static final String FILE_OUT_PATH = "
22
23      public static void main(String[] args) throws
24          URL resource = ConvertXmlToSequence.class
25
26          Configuration conf = new Configuration();
27
28          File inputFile = new File(resource.getPath
29          Path outputFile = new Path(resource.getPa
30
31          write(conf, inputFile, outputFile); // wr
32          read(conf, outputFile); // read seq file
33      }
34
35      /**
36       * Write a text file to sequence file
37       *
38       * @param conf
39       * @param inputFile
40       * @param outputFile
41       */
42      public static void write(Configuration conf,
43          SequenceFile.Writer writer = null;
44
45          try {
46              writer = SequenceFile.createWriter(co
47                      SequenceFile.Writer.compressi
48                      SequenceFile.Writer.keyClass(
49                      SequenceFile.Writer.valueClas
50
51              List<String> lines = FileUtils.readLi
52              String xmlString = lines.stream().map(
53              IntWritable key = new IntWritable(1);
54              BytesWritable value = new BytesWritab
55              writer.append(key, value);
56
57          } catch (IOException e) {
58              System.out.println("Error writing: "
59          }
60
```

```
61              finally {
62                  IOUtils.closeStream(writer);
63              }
64          }
65
66      /**
67       * Read a sequence file
68       *
69       * @param conf
70       * @param sequenceFileToRead
71       */
72      public static void read(Configuration conf, Pi
73          SequenceFile.Reader reader = null;
74          try {
75              reader = new SequenceFile.Reader(conf
76              IntWritable keyRead = (IntWritable) Re
77              BytesWritable valueRead = (BytesWritab
78              while (reader.next(keyRead, valueRead
79                  System.out.println("key : " + keyF
80                  valueRead = (BytesWritable) Refle
81              }
82          } catch (IOException e) {
83              System.out.println("Cannot read sequn
84          }
85
86          IOUtils.closeStream(reader);
87      }
88 }
89
```

Output:

```
1
2 key : 1 - value : <?xml version="1.0" encoding="UTI
3
```
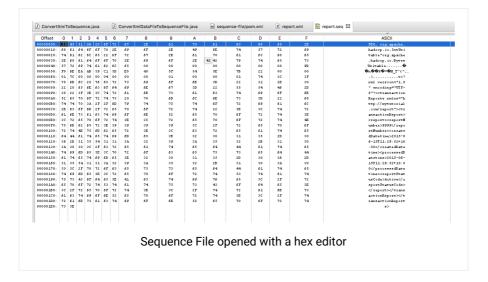
As you can see, the sequence files are stored as key/value pairs. The **key** being IntWritable "1" and the **value** being the "whole XML".

The **report.seq** file gets written to the project folder's **"target/classes/data"** folder. The "value" will be stored as binary.

Step 6: If you want to open the sequence file in eclipse, you need to install the hex editor plugin from

the update site:
"http://ehep.sourceforge.net/update".



Sequence File opened with a hex editor

if you write to HDFS, you can display it on HDFS shell with:

```
$ hdfs dfs -text path/report.seq
```

‹  02: Q7 – Q15 Hadoop overview & architecture interview Q&As

02: Convert XML file To Sequence File with Apache Spark – writing &

reading   ›

# Disclaimer