# 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 ▼

# 03: Convert XML file To an Avro File – writing & reading

📅 Posted on May 5, 2016

This extends the Convert XML file To Sequence File With Hadoop libaries. Avro files are **schema driven** & support **schema evolution**, which means you can add new columns & modify existing columns.

Step 1: The pom.xml file should include the Apache Spark libraries as shown below.

```
1  <project xmlns="http://maven.apache.org/POM/4.0.0"
2      xsi:schemaLocation="http://maven.apache.org/P(
3      <modelVersion>4.0.0</modelVersion>
4      <groupId>com.mytutorial</groupId>
5      <artifactId>sequence-file</artifactId>
```

## 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

```
 6        <packaging>jar</packaging>
 7        <version>1.0-SNAPSHOT</version>
 8        <name>sequence-file</name>
 9        <url>http://maven.apache.org</url>
10
11        <properties>
12            <maven.compiler.source>1.8</maven.compiler
13            <maven.compiler.target>1.8</maven.compiler
14            <project.build.sourceEncoding>UTF-8</proje
15            <junit.version>4.8.1</junit.version>
16            <hadoop.version>2.7.2</hadoop.version>
17            <spark-version>1.3.0</spark-version>
18        </properties>
19
20        <dependencies>
21            <!-- JUnit -->
22
23            <dependency>
24                <groupId>junit</groupId>
25                <artifactId>junit</artifactId>
26                <version>${junit.version}</version>
27                <scope>test</scope>
28            </dependency>
29
30            <!-- Hadoop -->
31            <dependency>
32                <groupId>org.apache.hadoop</groupId>
33                <artifactId>hadoop-hdfs</artifactId>
34                <version>${hadoop.version}</version>
35                <exclusions>
36                    <exclusion>
37                        <groupId>javax.servlet</group]
38                        <artifactId>*</artifactId>
39                    </exclusion>
40                </exclusions>
41            </dependency>
42            <dependency>
43                <groupId>org.apache.hadoop</groupId>
44                <artifactId>hadoop-client</artifactId>
45                <version>${hadoop.version}</version>
46                <exclusions>
47                    <exclusion>
48                        <groupId>javax.servlet</group]
49                        <artifactId>*</artifactId>
50                    </exclusion>
51                </exclusions>
52            </dependency>
53
54            <!-- Apache Spark -->
55            <dependency>
56                <groupId>org.apache.spark</groupId>
57                <artifactId>spark-core_2.11</artifactI
58                <version>${spark-version}</version>
59                <exclusions>
60                    <exclusion>
```

## 800+ Java Interview Q&As

```
61                        <groupId>javax.servlet</group]
62                        <artifactId>*</artifactId>
63                    </exclusion>
64                </exclusions>
65            </dependency>
66        </dependencies>
67
68  </project>
69
```

**Step 2:** The **report.xml** file under
"src/main/resources/data".

```
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+1
8               <processedDatetime>2015-06-15T11:29:52
9               <reportStatusCode>Active</reportStatus
10          </report>
11      </transactionReport>
12  </transactionReports>
13
```

**Step 3:** The avro schema file "**trans-report.avsc**" under
"src/main/resources/schema".

```
1
2   {"namespace": "mytutorial.com.report",
3    "type": "record",
4    "name": "ReportAvro",
5    "fields": [
6        {"name": "reportNumber", "type": "string"},
7        {"name": "createdDatetime", "type": "string"
8        {"name": "processedDatetime", "type": "string
9        {"name": "reportStatusCode", "type": "string"
10   ]
11  }
12
```

**Step 4:** The **Report.java** to map XML contents to
POJO (Plain Old Java Object).

```java
1
2    package com.mytutorial.pojo;
3
4    public class Report {
5
6        private String reportNumber;
7        private String createdDatetime;
8        private String processedDatetime;
9        private String reportStatusCode;
10
11       public Report(String reportNumber, String cre
12           this.reportNumber = reportNumber;
13           this.createdDatetime = createdDatetime;
14           this.processedDatetime = processedDatetime
15           this.reportStatusCode = reportStatusCode;
16       }
17
18       public String getReportNumber() {
19           return reportNumber;
20       }
21
22       public void setReportNumber(String reportNumbe
23           this.reportNumber = reportNumber;
24       }
25
26       public String getCreatedDatetime() {
27           return createdDatetime;
28       }
29
30       public void setCreatedDatetime(String createdI
31           this.createdDatetime = createdDatetime;
32       }
33
34       public String getProcessedDatetime() {
35           return processedDatetime;
36       }
37
38       public void setProcessedDatetime(String proces
39           this.processedDatetime = processedDatetime
40       }
41
42       public String getReportStatusCode() {
43           return reportStatusCode;
44       }
45
46       public void setReportStatusCode(String reportS
47           this.reportStatusCode = reportStatusCode;
48       }
49
50       @Override
51       public String toString() {
52           return "Report [reportNumber=" + reportNum
53                   + processedDatetime + ", reportSto
54       }
```

```
55 }
56
```

Step 5: Finally, the stand-alone
"ConvertXmlToAvroFile.java" to convert an XML to
POJO, and then to AVRO "GenericRecord", and then
to an AVRO file "data/report.avro".

```java
1
2    package com.mytutorial;
3
4    import java.io.File;
5    import java.io.IOException;
6    import java.io.StringReader;
7    import java.net.URL;
8    import java.util.Iterator;
9
10   import javax.xml.namespace.NamespaceContext;
11   import javax.xml.xpath.XPath;
12   import javax.xml.xpath.XPathConstants;
13   import javax.xml.xpath.XPathExpressionException;
14   import javax.xml.xpath.XPathFactory;
15
16   import org.apache.avro.Schema;
17   import org.apache.avro.file.DataFileReader;
18   import org.apache.avro.file.DataFileWriter;
19   import org.apache.avro.generic.GenericData;
20   import org.apache.avro.generic.GenericDatumReader
21   import org.apache.avro.generic.GenericDatumWriter
22   import org.apache.avro.generic.GenericRecord;
23   import org.apache.avro.io.DatumReader;
24   import org.apache.avro.io.DatumWriter;
25   import org.apache.commons.io.FileUtils;
26   import org.w3c.dom.Node;
27   import org.xml.sax.InputSource;
28
29   import com.mytutorial.pojo.Report;
30
31   public class ConvertXmlToAvroFile {
32
33       private static final String FILE_IN_PATH = "
34       private static final String FILE_OUT_PATH =
35       private static final String AVRO_SCHEMA_FILE
36
37       public static void main(String[] args) throws
38           URL resource = ConvertXmlToSequence.class
39
40           File inputFile = new File(resource.getPat
41           File outputFile = new File(resource.getP
42           File avroSchemaFile = new File(resource.
```

```
43
44          Report report = convertXmlToPojo(inputFi
45
46          write(avroSchemaFile, outputFile, report
47
48          read(avroSchemaFile, outputFile); // rea
49      }
50
51      // write the pojo "Report" to avro file
52      public static void write(File avroSchemaFile
53          Schema avroSchema = new Schema.Parser().
54          GenericRecord myrecord = new GenericData
55          myrecord.put("reportNumber", report.getR
56          myrecord.put("createdDatetime", report.g
57          myrecord.put("processedDatetime", report
58          myrecord.put("reportStatusCode", report.
59
60          DatumWriter<GenericRecord> datumWriter =
61          DataFileWriter<GenericRecord> writer = n
62          writer.create(avroSchema, outputFile);
63          writer.append(myrecord);
64          writer.close();
65      }
66
67      //read an avro file
68      private static void read(File avroSchemaFile
69          Schema avroSchema = new Schema.Parser().
70
71          DatumReader<GenericRecord> datumReader =
72          DataFileReader<GenericRecord> reader = n
73          GenericRecord record = reader.next();
74          System.out.println(record);
75          reader.close();
76      }
77
78      // Xpath to read XML & convert it to a pojo
79      private static Report convertXmlToPojo(File
80          XPath xPath = XPathFactory.newInstance()
81          // namespace
82          NamespaceContext ctx = new NamespaceCont
83              public String getNamespaceURI(String
84                  return prefix.equals("urn") ? "h
85              }
86
87              public Iterator<String> getPrefixes(
88                  return null;
89              }
90
91              public String getPrefix(String uri)
92                  return null;
93              }
94          };
95
96          xPath.setNamespaceContext(ctx);
97          String str = FileUtils.readFileToString(
```

```
 98         StringReader sr = new StringReader(str);
 99         InputSource source = new InputSource(sr)
100
101         // get the DOM
102         Node root = (Node) xPath.evaluate("/", s
103
104         // use the DOM
105         String reportNumber = xPath.evaluate("//
106         String createdDatetime = xPath.evaluate(
107         String processedDatetime = xPath.evaluat
108         String reportStatusCode = xPath.evaluate
109
110         return new Report(reportNumber, createdD
111     }
112 }
113
114
```

‹   XML Parsing with JAXB implementation called MOXy

    03: Q16 – Q26 Hadoop MapReduce interview questions & answers   ›

# Disclaimer

© 2022 java-success.com