# 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 ▾

# 2. Apache Pig: Regex (Regular expressions)

📅 Posted on February 5, 2016

This extends the tutorial 1. Apache Pig Getting started.

## Input Data

**scores.xml** in folder:**/Users/arulk/projects** representing marks of 4 students in 3 subjects:

```
1
2   <scores>
3     <subject>
4       <name>Science</name>
5       <marks>
6         <mark>80</mark>
7       <mark>75</mark>
```

## 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

```
 8          <mark>89</mark>
 9          <mark>90</mark>
10        </marks>
11      </subject>
12      <subject>
13          <name>Maths</name>
14        <marks>
15            <mark>90</mark>
16        <mark>87</mark>
17        <mark>78</mark>
18        <mark>92</mark>
19        </marks>
20      </subject>
21      <subject>
22          <name>English</name>
23        <marks>
24            <mark>78</mark>
25        <mark>88</mark>
26        <mark>65</mark>
27        <mark>99</mark>
28        </marks>
29      </subject>
30  </scores>
31
```

**Step 1:** Start pig in local file system mode.

```
1
2 pig -x local
3
```

**Step 2:** Extract the "**Subjects**" from the input XML file.

```
1
2 grunt> SUBJECTS_EXTRACT =  LOAD '/Users/arulk/proj
3
```

Dump the output:

```
1
2 grunt> dump SUBJECTS_EXTRACT;
3
```

```
1
2 (<subject>        <name>Science</name>      <marks
3 (<subject>        <name>Maths</name>        <marks
```

```
4  (<subject>         <name>English</name>        <marks
5
```

## Step 3: Regex to extract each "Subject" and its corresponding marks.

```
1
2  grunt> MARKS_FOR_SUBJECT_CSV = foreach SUBJECTS_EX
3
```

```
1
2  grunt> dump MARKS_FOR_SUBJECT_CSV;
3
```

Dump it:

```
1
2  (Science,80,75,89,90)
3  (Maths,90,87,78,92)
4  (English,78,88,65,99)
5
```

## Step 4: Put it all into a single "marks_by_subjects.pig" script.

```
1
2
3  SUBJECTS_EXTRACT =  LOAD '/Users/arulk/projects/sc
4
5  MARKS_FOR_SUBJECT_CSV = foreach SUBJECTS_EXTRACT   (
6
7  dump MARKS_FOR_SUBJECT_CSV
8
9
```

Run the above pig script:

```
1
2  $ pig -x local marks_by_subjects.pig
3
```

## Outputs:

```
1
2  (Science,80,75,89,90)
3  (Maths,90,87,78,92)
4  (English,78,88,65,99)
5
```

If you run without "-x local" option it runs in "Map-Reduce" mode against the HDFS (e.g. hdfs://localhost:9000). The name and data nodes need to be running. Otherwise you will get "Caused by: java.net.ConnectException: Connection refused" error.

## Mapreduce mode

```
1
2  pig -x mapreduce
3  pig
4
```

‹   1. Apache Pig Getting started        3. Apache Pig: XPath for XML   ›

## Disclaimer

© 2022 java-success.com

Top