%spark

```
%dep

z.reset()
z.load("joda-time:joda-time:2.9.1")
```

DepInterpreter(%dep) deprecated. Remove dependencies and repositories through GUI interpreter menu instead.

DepInterpreter(%dep) deprecated. Load dependency through GUI interpreter menu instead. res0: org.apache.zeppelin.dep.Dependency = org.apache.zeppelin.dep.Dependency@3b8b1761

Took 7 sec. Last updated by anonymous at February 02 2017, 8:29:42 PM.

FINISHED ▷ 光 圓 ��

```
import org.apache.spark.rdd._
import scala.collection.JavaConverters._
import au.com.bytecode.opencsv.CSVReader
import java.sql.Timestamp
import java.io._
import org.joda.time._
import org.joda.time.format._
import org.joda.time.format.DateTimeFormat
import org.joda.time.DateTime
import org.joda.time.Days
case class DelayRec(year: String,
                       month: String,
                       dayOfMonth: String,
                       dayOfWeek: String,
                       crsDepTime: String,
                       depDelay: String,
                       origin: String,
                       distance: String,
                       cancelled: String) {
    val holidays = List("01/01/2007", "01/15/2007", "02/19/2007", "05/28/2007", "06/07/2007",
       "09/03/2007", "10/08/2007", "11/11/2007", "11/22/2007", "12/25/2007", "01/01/2008", "01/21/2008", "02/18/2008", "05/22/2008", "05/26/2008", "07/04/2008", "09/01/2008", "10/13/2008", "11/11/2008", "11/27/2008", "12/25/2008")
    def gen_features: (String, Array[Double]) = {
       val values = Array(
         depDelay.toDouble,
         month.toDouble,
         dayOfMonth.toDouble,
         dayOfWeek.toDouble,
         get_hour(crsDepTime).toDouble,
         distance.toDouble,
```

```
days_from_nearest_holiday(year.toInt, month.toInt, dayOfMonth.toInt)
       )
       new Tuple2(to_date(year.toInt, month.toInt, dayOfMonth.toInt), values)
     def get_hour(depTime: String) : String = "%04d".format(depTime.toInt).take(2)
     def to_date(year: Int, month: Int, day: Int) = "%04d%02d%02d".format(year, month, day)
     def days_from_nearest_holiday(year:Int, month:Int, day:Int): Int = {
       val sampleDate = new org.joda.time.DateTime(year, month, day, 0,0)
       holidays.foldLeft(3000) \{ (r, c) =>
         val holiday = org.joda.time.format.DateTimeFormat.forPattern("MM/dd/yyyy").parseDate
         val distance = Math.abs( org.joda.time.Days.daysBetween(holiday, sampleDate).getDays)
         math.min(r, distance)
       }
     }
   }
import org.apache.spark.rdd._
import scala.collection.JavaConverters._
import au.com.bytecode.opencsv.CSVReader
import java.sql.Timestamp
import java.io._
import org.joda.time._
import org.joda.time.format._
import org.joda.time.format.DateTimeFormat
import org.joda.time.DateTime
import org.joda.time.Days
defined class DelayRec
Took 8 sec. Last updated by anonymous at February 02 2017, 8:48:13 PM.
```

```
FINISHED ▷ 💥 🗐 🅸
                  // function to do a preprocessing step for a given file
                   def prepFlightDelays(infile: String): RDD[DelayRec] = {
                                  val data = sc.textFile(infile)
                                  data.map { line =>
                                         val reader = new CSVReader(new StringReader(line))
                                          reader.readAll().asScala.toList.map(rec => DelayRec(rec(0),rec(1),rec(2),rec(3),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5),rec(5)
                                  }.map(list => list(0))
                                  .filter(rec => rec.year != "Year")
                                   .filter(rec => rec.cancelled == "0")
                                   .filter(rec => rec.origin == "ORD")
                   }
Projecta_207tmp = prepFlightDelays("/Users/vpandiyan/Downloads/flights_2007.csv.bz2")
                    val data_2008 = prepFlightDelays("/Users/vpandiyan/Downloadss/flights_2008.csv.bz2").map(rec =
                    data_2007tmp.toDF().registerTempTable("data_2007tmp")
                   data_2007.take(5).map(x \Rightarrow x mkString ",").foreach(println)
```

prepFlightDelays: (infile: String)org.apache.spark.rdd.RDD[DelayRec]

data_2007tmp: org.apache.spark.rdd.RDD[DelayRec] = MapPartitionsRDD[57] at filter at <console</pre>

>:119

data_2007: org.apache.spark.rdd.RDD[Array[Double]] = MapPartitionsRDD[58] at map at <console>:

data_2008: org.apache.spark.rdd.RDD[Array[Double]] = MapPartitionsRDD[66] at map at <console>:
111

warning: there was one deprecation warning; re-run with -deprecation for details

-8.0,1.0,25.0,4.0,11.0,719.0,10.0

41.0,1.0,28.0,7.0,15.0,925.0,13.0

45.0,1.0,29.0,1.0,20.0,316.0,14.0

-9.0,1.0,17.0,3.0,19.0,719.0,2.0

180.0,1.0,12.0,5.0,17.0,316.0,3.0

Took 9 sec. Last updated by anonymous at February 02 2017, 8:48:26 PM.

%sql

select * from data_2007tmp



FINISHED ▷ ତ 🕸

year	month	dayOfMonth	dayOfWeek	crsDepTime	depDelay
2,007	1	19	5	1,210	46
2,007	1	11	4	1,500	-7
2,007	1	23	2	1,500	-4
2,007	1	3	3	1,210	-9
2,007	1	12	5	1,330	103
2,007	1	8	1	1,255	16
2,007	1	26	5	1,100	-5
2,007	1	30	2	1,210	27
2,007	1	19	5	602	-3

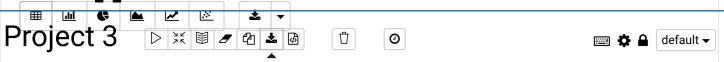
Results are limited by 1000.

Took 3 sec. Last updated by anonymous at February 02 2017, 8:48:38 PM.

%sql FINISHED ▷ 光 ፱ ⑫

Profest_day fWeek, case when depDelay > 15 then 'delayed' else 'ok' end , count(1)

up respection case when depDelay > 15 then 'delayed' else 'ok' end



Export the notebook

dayofWeek	CASE WHEN (CAST(depDelay AS DOUBLE) > CAST(15 AS DOUBLE)) THEN delayed EL
1	delayed
7	ok
1	ok
6	delayed
2	delayed
3	ok
4	delayed
3	delayed
5	ok

%sql

FINISHED ▷ 💥 🗊 🕸

select cast(cast(crsDepTime as int) / 100 as int) as hour, case when depDelay > 15 then 'del count

from data_2007tmp

group by cast(cast(crsDepTime as int) / 100 as int), case when depDelay > 15 then 'delayed'



hour	delay	c
12	ok	1
13	ok	2
20	delayed	1
10	ok	1
19	ok	1
15	ok	1
15	delayed	7
21	ok	8
8	ok	2

mi**zeppelim**onymous at February 02 2017, 8:52:16 PM.

