

```
%dep
```

FINISHED ▶ ⌵ 📖 ⚙

```
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@7d520ec8

Took 7 sec. Last updated by anonymous at February 02 2017, 10:38:53 PM.

```
// Imports
import org.apache.spark.sql.functions._
import org.joda.time.format.DateTimeFormat
```

FINISHED ▶ ⌵ 📖 ⚙

```
import org.apache.spark.sql.functions._
import org.joda.time.format.DateTimeFormat
```

Took 3 sec. Last updated by anonymous at February 02 2017, 10:39:06 PM.

```
// Load data - adjust the path to the location of your data
```

FINISHED ▶ ⌵ 📖 ⚙

```
val inputPath_2007 = "/Users/neha/Documents/Capstone/Week1-zeppelin_getting_started/Data/year/2007.csv"
val inputPath_2008 = "/Users/neha/Documents/Capstone/Week1-zeppelin_getting_started/Data/year/2008.csv"
val inputPath_we_2007 = "/Users/neha/Documents/Capstone/Week3/data/2007.csv"
val inputPath_we_2008 = "/Users/neha/Documents/Capstone/Week3/data/2008.csv"
```

inputPath_2007: String = /Users/neha/Documents/Capstone/Week1-zeppelin_getting_started/Data/year/2007.csv

inputPath_2008: String = /Users/neha/Documents/Capstone/Week1-zeppelin_getting_started/Data/year/2008.csv

inputPath_we_2007: String = /Users/neha/Documents/Capstone/Week3/data/2007.csv

inputPath_we_2008: String = /Users/neha/Documents/Capstone/Week3/data/2008.csv

Took 2 sec. Last updated by anonymous at February 02 2017, 10:39:10 PM.

```
% spark
import org.apache.spark.rdd._
import scala.collection.JavaConverters._
import au.com.bytecode.opencsv.CSVReader
```

FINISHED ▶ ⌵ 📖 ⚙

```
res4: org.apache.spark.sql.Session = org.apache.spark.sql.Session@71a09c97
import org.apache.spark.rdd._
import scala.collection.JavaConverters._
```

```
import au.com.bytecode.opencsv.CSVReader
```

Took 2 sec. Last updated by anonymous at February 02 2017, 10:39:14 PM.

FINISHED ▶ ⌕ 📖 ⚙

```
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
import org.joda.time.format.DateTimeFormat
import org.joda.time.format.DateTimeFormatter
```

```
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
import org.joda.time.format.DateTimeFormat
import org.joda.time.format.DateTimeFormatter
```

Flight delay analysis-Week3

Took 2 sec. Last updated by anonymous at February 02 2017, 10:39:19 PM. | Untitled | Untitled | Untitled | Untitled | Untitled | Untitled | Untitled |

FINISHED ▶ ⌕ 📖 ⚙

```
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").parseDateT
  val distance = Math.abs(org.joda.time.Days.daysBetween(holiday, sampleDate).getDays)
  math.min(r, distance)
}
}
}

```

defined class DelayRec

Took 2 sec. Last updated by anonymous at February 02 2017, 10:39:24 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),1
    ).map(rec => DelayRec(rec(0),rec(1),rec(2),rec(3),rec(5),1
    ).filter(rec => rec.year != "Year")
    .filter(rec => rec.cancelled == "0")
    .filter(rec => rec.origin == "ORD")
  }
}

```

Flight delay analysis-Week3

Zeppelin

```

val data_2007tmp = prepFlightDelays(inputPath_2007)
val data_2007 = data_2007tmp.map(rec => rec.gen_features._2)
val data_2008 = prepFlightDelays(inputPath_2008).map(rec => rec.gen_features._2)

```

Flight delay analysis...

```
data_2007tmp.toDF().registerTempTable("data_2007tmp")
```

```
data_2007.take(5).map(x => x.mkString ",").foreach(println)
```

```

prepFlightDelays: (infile: String)org.apache.spark.rdd.RDD[DelayRec]
data_2007tmp: org.apache.spark.rdd.RDD[DelayRec] = MapPartitionsRDD[68] at filter at <console>:64
data_2007: org.apache.spark.rdd.RDD[Array[Double]] = MapPartitionsRDD[69] at map at <console>:60
data_2008: org.apache.spark.rdd.RDD[Array[Double]] = MapPartitionsRDD[77] at map at <console>:58
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 4 sec. Last updated by anonymous at February 02 2017, 10:55:18 PM.

FINISHED ▶ ⌕ 📖 ⚙️

```
%sql
```

```
select dayofWeek, case when depDelay > 15 then 'delayed' else 'ok' end , count(1)
```

```
from data_2007tmp group by dayofweek case when depDelay > 15 then 'delayed' else 'ok' end
```





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

Flight delay analysis-Week3

```
%[select cast(cast(crsDepTime as int) / 100 as int) as hour, case when depDelay > 15 then 'del'
from data_2007tmp group by cast( cast(crsDepTime as int) / 100 as int), case when depDelay
```

FINISHED ▶ ⌵ 📖 ⚙️

Flight delay analysis





default ▼

hour	delay
12	ok
13	ok
20	delayed
10	ok
19	ok
15	ok
15	delayed
21	ok
8	ok

Took 23 sec. Last updated by anonymous at February 02 2017, 10:58:23 PM.

READY ▶ ⌵ 📖 ⚙️

Flight delay analysis-Week3



Zeppelin

Flight delay analysis...

default