

2017-02-02-Flight-D...



default ▾

1 %sh

2

3 wget http://stat-computing.org/dataexpo/2009/2007.csv.bz2 -O /Users/Hamster/Desktop/Air_F

4 wget http://stat-computing.org/dataexpo/2009/2008.csv.bz2 -O /Users/Hamster/Desktop/Air_F

5 echo "download"

FINISHED ▶ ⌵ 📖 ⚙

--2017-02-02 23:20:28-- http://stat-computing.org/dataexpo/2009/2007.csv.bz2 ↴

Resolving stat-computing.org (stat-computing.org)... 52.218.128.147

Connecting to stat-computing.org (stat-computing.org)|52.218.128.147|:80... connected.

HTTP request sent, awaiting response... 200 OK

Length: 121249243 (116M) [application/x-bzip2]

Took 1 min 1 sec. Last updated by anonymous at February 02 2017, 11:21:28 PM.

1 %sh

2

3 wget ftp://ftp.ncdc.noaa.gov/pub/data/ghcn/daily/by_year/2007.csv.gz -O /Users/Hamster/De

4 wget ftp://ftp.ncdc.noaa.gov/pub/data/ghcn/daily/by_year/2008.csv.gz -O /Users/Hamster/De

5 echo "download"

FINISHED ▶ ⌵ 📖 ⚙

--2017-02-02 23:21:51-- ftp://ftp.ncdc.noaa.gov/pub/data/ghcn/daily/by_year/2007.csv.gz ↴

=> '/Users/Hamster/Desktop/Air_Flight_Data/weather_2007.csv.gz'

Resolving ftp.ncdc.noaa.gov (ftp.ncdc.noaa.gov)... 205.167.25.101, 2610:20:8040:2::101

Connecting to ftp.ncdc.noaa.gov (ftp.ncdc.noaa.gov)|205.167.25.101|:21... connected.

Logging in as anonymous ... Logged in!

==> SYST ... done. ==> PWD ... done.

==> TYPE I ... done. ==> CWD (1) /pub/data/ghcn/daily/by_year ... done.

==> SIZE 2007.csv.gz ... 197708335

==> PASV ... done. ==> RETR 2007.csv.gz ... done.

Length: 197708335 (189M) (unauthoritative)

0K	0%	674K	4m46s
50K	0%	907K	4m9s
100K	0%	1.35M	3m33s
150K	0%	1.34M	3m15s
200K	0%	1.38M	3m3s
250K	0%	12.4M	2m35s
300K	0%	1.79M	2m28s
350K	0%	1.55M	2m21s

Took 1 min 0 sec. Last updated by anonymous at February 02 2017, 11:22:51 PM.

1 %dep

2

3 z.reset()

4 z.load("joda-time:joda-time:2.9.1")

FINISHED ▶ ⌵ 📖 ⚙

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@4eb3722a

Zeppelin

2017-02-Flight-D...

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

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

Took 5 sec. Last updated by anonymous at February 02 2017, 11:06:27 PM.

FINISHED ▶ 🔍 📖 ⚙️
 🖨️ ⚙️ 🔒 default ▼

```
1 import java.io._
2 import org.joda.time._
3 import org.joda.time.format._
4 import org.joda.time.format.DateTimeFormat
5 import org.joda.time.DateTime
6 import org.joda.time.Days
```

```
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
```

Took 1 sec. Last updated by anonymous at February 02 2017, 11:06:31 PM.

FINISHED ▶ ⌵ ⌶ ⚙

```

1 case class DelayRec(year: String,
2                       month: String,
3                       dayOfMonth: String,
4                       dayOfWeek: String,
5                       crsDepTime: String,
6                       depDelay: String,
7                       origin: String,
8                       distance: String,
9                       cancelled: String) {
10
11   val holidays = List("01/01/2007", "01/15/2007", "02/19/2007", "05/28/2007", "06/07/2007",
12                      "09/03/2007", "10/08/2007", "11/11/2007", "11/22/2007", "12/25/2007",
13                      "01/01/2008", "01/21/2008", "02/18/2008", "05/22/2008", "05/26/2008", "07/04/2008",
14                      "09/01/2008", "10/13/2008", "11/11/2008", "11/27/2008", "12/25/2008")
15
16   def gen_features: (String, Array[Double]) = {
17     val values = Array(
18       depDelay.toDouble,
19       month.toDouble,
20       dayOfMonth.toDouble,
21       dayOfWeek.toDouble,
22       get_hour(crsDepTime).toDouble,
23       distance.toDouble,
24       days_from_nearest_holiday(year.toInt, month.toInt, dayOfMonth.toInt)
25     )
26     new Tuple2(to_date(year.toInt, month.toInt, dayOfMonth.toInt), values)

```

FINISHED ▶ ⌵ ⌶ ⚙

Zeppelin

2017-02-02-Flight-D

```

29 def get_hour(depTime: String) : String = "%04d".format(depTime.toInt).take(2)
30 def to_date(year: Int, month: Int, day: Int) = "%04d%02d%02d".format(year, month, dc
31
32 def days_from_nearest_holiday(year: Int, month: Int, day: Int): Int = {
33     val sampleDate = new org.joda.time.DateTime(year, month, day, 0, 0)
34
35     holidays.foldLeft(3000) { (r, c) =>
36         val holiday = org.joda.time.format.DateTimeFormat.forPattern("MM/dd/yyyy").parse
37         val distance = Math.abs(org.joda.time.Days.daysBetween(holiday, sampleDate).getL
38         math.min(r, distance)
39     }
40 }
```

defined class DelayRec

Took 1 sec. Last updated by anonymous at February 02 2017, 11:06:34 PM.

FINISHED ▶ ⌕ 📖 ⚙

```

1 // function to do a preprocessing step for a given file
2 def prepFlightDelays(infile: String): RDD[DelayRec] = {
3     val data = sc.textFile(infile)
4
5     data.map { line =>
6         val reader = new CSVReader(new StringReader(line))
7         reader.readAll().asScala.toList.map(rec => DelayRec(rec(0),rec(1),rec(2),rec(3),re
8     }.map(list => list(0))
9     .filter(rec => rec.year != "Year")
10    .filter(rec => rec.cancelled == "0")
11    .filter(rec => rec.origin == "ORD")
12 }
```

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

Took 0 sec. Last updated by anonymous at February 02 2017, 11:06:37 PM.

FINISHED ▶ ⌕ 📖 ⚙

```

1 val data_2007tmp = prepFlightDelays("/Users/Hamster/Desktop/Air_Flight_Data/flights_2007.
2 val data_2007 = data_2007tmp.map(rec => rec.gen_features._2)
3 val data_2008 = prepFlightDelays("/Users/Hamster/Desktop/Air_Flight_Data/flights_2008.csv
4
5 data_2007tmp.toDF().registerTempTable("data_2007tmp")
6
7 data_2007.take(5).map(x => x mkString ",").foreach(println)
```

data_2007tmp: org.apache.spark.rdd.RDD[DelayRec] = MapPartitionsRDD[6] at filter at <console>:

58

data_2007: org.apache.spark.rdd.RDD[Array[Double]] = MapPartitionsRDD[7] at map at <console>:5

2

data_2008: org.apache.spark.rdd.RDD[Array[Double]] = MapPartitionsRDD[15] at map at <console>:

50

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 8 sec. Last updated by anonymous at February 02 2017, 11:06:49 PM.

2017-02-02-Flight-Data

Zeppelin

1 %sql
2 select dayofWeek, case when depDelay > 15 then 'delayed' else 'ok' end , count(1)
3 from data_2007tmp group by dayofweek , case when depDelay > 15 then 'delayed' else 'ok' end

FINISHED ▶ ⌵ 📖 ⚙

▶ ⌵ 📖 ✎ 📄 ⬇️ 📄 🗑️ ⌚ ⌨️ ⚙ 🔒 default ▼

📊 📈 📉 📱 📈 📊 ⬇️ ▼

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

Took 25 sec. Last updated by anonymous at February 02 2017, 11:07:32 PM.

1 %sql
2 select cast(cast(crsDepTime as int) / 100 as int) as hour, case when depDelay > 15 then 'delayed' else 'ok' end

FINISHED ▶ ⌵ 📖 ⚙

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, 11:09:18 PM.

2017-02-02-Flight-Data

1 %spark

FINISHED ▶ ⌵ 📖 ⚙

Zeppelin

```
import org.apache.spark.mllib.linalg.Vectors
```

```
import org.apache.spark.mllib.feature.StandardScaler
```

2017-02-02-Flight-D.

01/02/02 - Night 0

```
import org.apache.spark.mllib.linalg.Vectors
```

```
import org.apache.spark.mllib.feature.StandardScaler
```

Took 0 sec. Last updated by anonymous at February 02 2017, 11:09:25 PM.

```
1 def parseData(vals: Array[Double]): LabeledPoint = {
2     LabeledPoint(if (vals(0))>=15) 1.0 else 0.0, Vectors.dense(vals.drop(1)))
3 }
```

```
parseData: (vals: Array[Double])org.apache.spark.mllib.regression.LabeledPoint
```

Took 0 sec. Last updated by anonymous at February 02 2017, 11:09:28 PM.

```
1 // Prepare training set
2 val parsedTrainData = data_2007.map(parseData)
3 parsedTrainData.cache
4 val scaler = new StandardScaler(withMean = true, withStd = true).fit(parsedTrainData.map(
5 val scaledTrainData = parsedTrainData.map(x => LabeledPoint(x.label, scaler.transform(Vec
6 scaledTrainData.cache
```

```
parsedTrainData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = MapPartitionsRDD[48] at map at <console>:60
```

```
res11: parsedTrainData.type = MapPartitionsRDD[48] at map at <console>:60
```

```
scaler: org.apache.spark.mllib.feature.StandardScalerModel = org.apache.spark.mllib.feature.St
andardScalerModel@4b73aaf8
```

```
scaledTrainData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = MapPartitionsRDD[51] at map at <console>:63
```

```
res12: scaledTrainData.type = MapPartitionsRDD[51] at map at <console>:63
```

Took 28 sec. Last updated by anonymous at February 02 2017, 11:09:58 PM.

```
1 // Prepare test/validation set
2 val parsedTestData = data_2008.map(parseData)
3 parsedTestData.cache
4 val scaledTestData = parsedTestData.map(x => LabeledPoint(x.label, scaler.transform(Vector(x.features))))
5 scaledTestData.cache
```

```
parsedTestData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = Map
PartitionsRDD[52] at map at <console>:58
```

```
res15: parsedTestData.type = MapPartitionsRDD[52] at map at <console>:58
```

```
scaledTestData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = Map
PartitionsRDD[53] at map at <console>:67
```

```
res16: scaledTestData.type = MapPartitionsRDD[53] at map at <console>:67
```

Took 1 sec. Last updated by anonymous at February 02 2017, 11:10:12 PM.

```
1 scaledTrainData.take(3).map(x => (x.label, x.features)).foreach(println)
```

(0.0, [-1.6160463330366632, 1.0549272994666004, 0.03217026353736743, -0.518924417544128, 0.0340833342430724, -0.28016830994663705])

017.02.02 Flight Data

FINISHED

2017-02-02 Build the Logistic Regression Model

```

4 // Build the Logistic Regression model
5 val model_lr = LogisticRegressionWithSGD.train(scaledTrainData, numIterations=100)
6
7 // Predict
8 val labelsAndPreds_lr = scaledTestData.map { point =>
9     val pred = model_lr.predict(point.features)
10     (pred, point.label)
11 }
12
13 val m_lr = eval_metrics(labelsAndPreds_lr)._2

```

Took 30 sec. Last updated by anonymous at February 02 2017, 11:11:02 PM.

FINISHED ▶ ⌵ ⌶ ⚙

Took 0 sec. Last updated by anonymous at February 02 2017, 11:11:22 PM.

FINISHED ▶ ⌵ ⌶ ⚙

↓

2017-02-02-Flight-D...

Took 3 sec. Last updated by anonymous at February 02 2017, 11:11:31 PM.

FINISHED

Took 17 sec. Last updated by anonymous at February 02 2017, 11:12:03 PM.

FINISHED ▶ ⌵ ⌶ ⚙

Took 14 sec. Last updated by anonymous at February 02 2017, 11:12:21 PM.

FINISHED ▶ ↻ 📖 ⚙️

2017-02-Flight-D...

```

3 import org.apache.spark.mllib.regression.LabeledPoint
4 import org.apache.spark.mllib.linalg.Vectors
5 import org.apache.spark.mllib.feature.StandardScaler
6
7 def parseData(vals: Array[Double]): LabeledPoint = {
8   LabeledPoint(if (vals(0)>=15) 1.0 else 0.0, Vectors.dense(vals.drop(1)))
9 }
10
11 // Prepare training set
12 val parsedTrainData = data_2007.map(parseData)
13 val scaler = new StandardScaler(withMean = true, withStd = true).fit(parsedTrainData.map
14 val scaledTrainData = parsedTrainData.map(x => LabeledPoint(x.label, scaler.transform(Vect
15 parsedTrainData.cache
16 scaledTrainData.cache
17
18 // Prepare test/validation set
19 val parsedTestData = data_2008.map(parseData)
20 val scaledTestData = parsedTestData.map(x => LabeledPoint(x.label, scaler.transform(Vect
21 parsedTestData.cache
22 scaledTestData.cache
23
24 scaledTrainData.take(5).map(x => (x.label, x.features)).foreach(println)

```

```
import org.apache.spark.mllib.regression.LabeledPoint
import org.apache.spark.mllib.linalg.Vectors
import org.apache.spark.mllib.feature.StandardScaler
parseData: (vals: Array[Double])org.apache.spark.mllib.regression.LabeledPoint
parsedTrainData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = MapPartitionsRDD[320] at map at <console>:80
scaler: org.apache.spark.mllib.feature.StandardScalerModel = org.apache.spark.mllib.feature.StandardScalerModel@75704fdf
scaledTrainData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = MapPartitionsRDD[323] at map at <console>:84
res36: parsedTrainData.type = MapPartitionsRDD[320] at map at <console>:80
res37: scaledTrainData.type = MapPartitionsRDD[323] at map at <console>:84
parsedTestData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = MapPartitionsRDD[324] at map at <console>:80
scaledTestData: org.apache.spark.rdd.RDD[org.apache.spark.mllib.regression.LabeledPoint] = MapPartitionsRDD[325] at map at <console>:88
res38: parsedTestData.type = MapPartitionsRDD[324] at map at <console>:80
res39: scaledTestData.type = MapPartitionsRDD[325] at map at <console>:88
```

Took 4 sec. Last updated by anonymous at February 02 2017, 11:13:52 PM.

```
1 %spark
2
3 import org.apache.spark.mllib.classification.LogisticRegressionWithSGD
4
5 // Build the Logistic Regression model
6 val model_lr = LogisticRegressionWithSGD.train(scaledTrainData, numIterations=100)
7
8 // Predict
9 val labelsAndPreds_lr = scaledTestData.map { point =>
```

2017-02-02

```
import org.apache.spark.mllib.classification.LogisticRegressionWithSGD
warning: there was one deprecation warning; re-run with -deprecation for details
model_lr: org.apache.spark.mllib.classification.LogisticRegressionModel = org.apache.spark.mllib.classification.LogisticRegressionModel: intercept = 0.0, numFeatures = 11, numClasses = 2, threshold = 0.5
labelsAndPreds_lr: org.apache.spark.rdd.RDD[(Double, Double)] = MapPartitionsRDD[475] at map at t <console>:95
m_lr: Metrics = Metrics@452ddebdb
precision = 0.40, recall = 0.68, F1 = 0.50, accuracy = 0.62
```

Took 1 min 3 sec. Last updated by anonymous at February 02 2017, 11:15:00 PM.

```
1 %spark
2
3 println(model_lr.weights)
```

$[-0.002641910361922112, 0.016515122840474603, -0.021125675242088107, 0.42541032171918614, 0.04784674910826409, 0.010015678133278754, 0.02539098796205592, -0.15003098865096928, 0.2829929186564463, 0.2365535217050361, 0.15029474134967363]$

Took 1 sec. Last updated by anonymous at February 02 2017, 11:15:09 PM.

```
1 %spark
2
3 import org.apache.spark.mllib.tree.DecisionTree
4
5 // Build the Decision Tree model
6 val numClasses = 2
7 val categoricalFeaturesInfo = Map[Int, Int]()
8 val impurity = "gini"
9 val maxDepth = 10
10 val maxBins = 100
11 val model_dt = DecisionTree.trainClassifier(parsedTrainData, numClasses, categoricalFeat
12
13 // Predict
14 val labelsAndPreds_dt = parsedTestData.map { point =>
15     val pred = model_dt.predict(point.features)
16     (point.label, pred)
17 }
18 val m_dt = new Metrics(labelsAndPreds_dt)
19 println("precision = %.2f, recall = %.2f, F1 = %.2f, accuracy = %.2f"
20         .format(m_dt.precision, m_dt.recall, m_dt.F1, m_dt.accuracy))
```

```
import org.apache.spark.mllib.tree.DecisionTree
numClasses: Int = 2
categoricalFeaturesInfo: scala.collection.immutable.Map[Int,Int] = Map()
impurity: String = gini
maxDepth: Int = 10
maxBins: Int = 100
model_dt: org.apache.spark.mllib.tree.model.DecisionTreeModel = DecisionTreeModel classifier o
f depth 10 with 1861 nodes
```

Zeppelin

2017-02-02-Flight-D





2017-02-02-Flight-D

2017-02-02-Flight-D

FINISHED    

↓

Took 49 sec. Last updated by anonymous at February 02 2017, 11:16:48 PM.

READY    

↓