

Problem Set 6

Alexander Brandt

SID: 24092167

October 19 2015

Friendly Collaborators: Milos Atz, Alex Ojala. Most of our discussion centered around RSQLite and Spark syntax. Though I did convince some new folks to try CyberDuck.

1

The file takes almost 40 minutes to generate on an xl.large instance, and balloons from 1.7 Gb to about 18 Gb. The code to generate is detailed in the main chunk.

Listing 1: ls results

```
-rw-r--r--  1 ubuntu ubuntu  18G Nov  2 12:14 FlightDatabase.sqlite
```

2

2.1

The filtering step (detailed in the main code, below) is pretty much the same for both steps: it consists of removing NA's, and in the case of spark, getting rid of header lines from the CSV. Both are easily accomplished with a few lines to remove the offending files. In the case of R, we tag each NA as a numeric code (0.1234) before deleting from the database.

2.2

The code for Spark and PySpark is largely modified from what Chris gave us in Unit 7. Instead of computing a "median" for the time delay, we instead bin the times, then map the key/tuple to a string and write it to a file in the hadoop filesystem. By and large, the SPARK method is incredibly faster (all operations can be performed within a few minutes, vs. the arduously long loading times needed by R). For Spark, the output has been slightly modified to highlight the runtimes, without the verbosity.

For SPARK/PySPARK:

```
import time
from operator import add
import numpy as np
from pyspark import SparkContext

sc = SparkContext()

print "Read in start/stop times..."
print time.strftime("%H:%M:%S")
lines = sc.textFile('/data/airline')
print time.strftime("%H:%M:%S")
```



```
print time.strftime("%H:%M:%S")
myResults = output.map(binFun)
print time.strftime("%H:%M:%S")
myResults.map(printable).repartition(1).saveAsTextFile('/data/airline_processed')
```

Output:

Listing 2: ls results

```
15/11/02 10:15:47 INFO spark.SparkContext: Running Spark version 1.5.1
15/11/02 10:15:48 WARN spark.SparkConf:
SPARK_WORKER_INSTANCES was detected (set to '1').
This is deprecated in Spark 1.0+.
```

Please instead use:

- ./spark-submit with --num-executors to specify the number of executors
- Or `set SPARK_EXECUTOR_INSTANCES`
- `spark.executor.instances` to configure the number of instances in the spark config.

```
15/11/02 10:15:48 INFO spark.SecurityManager: Changing view acls to: root
15/11/02 10:15:48 INFO spark.SecurityManager: Changing modify acls to: root
15/11/02 10:15:48 INFO spark.SecurityManager: SecurityManager: authentication disabled
15/11/02 10:15:49 INFO slf4j.Slf4jLogger: Slf4jLogger started
15/11/02 10:15:49 INFO Remoting: Starting remoting
15/11/02 10:15:49 INFO Remoting: Remoting started; listening on addresses :[akka.tcp:/
15/11/02 10:15:49 INFO util.Utils: Successfully started service 'sparkDriver' on port
15/11/02 10:15:49 INFO spark.SparkEnv: Registering MapOutputTracker
15/11/02 10:15:49 INFO spark.SparkEnv: Registering BlockManagerMaster
[...]
15/11/02 10:26:13 INFO remote.RemoteActorRefProvider$RemotingTerminator: Remote daemon
```

And if you wanted to see a few lines from the built file:

Listing 3: ls results

```
Alexanders-MBP:airline_processed Alex$ head part-00000
DL,JFK,ORD,10,6,18,2,0,0,24,0.083333
YV,CLT,GSO,1,7,00,1,1,1,1,1.000000
AA,HDN,DFW,3,7,13,2,0,0,74,0.027027
OH,JFK,DTW,5,7,19,1,1,0,9,0.111111
OO,ICT,DEN,2,5,14,1,1,0,1,1.000000
AA,DFW,ATL,9,5,12,4,3,0,90,0.044444
EV,ATL,SHV,10,2,13,0,0,0,3,0.000000
US,SFO,PHL,1,6,22,0,0,0,17,0.000000
HP,LAS,MIA,1,3,22,0,0,0,1,0.000000
AS,SEA,GEG,9,5,23,2,1,0,32,0.062500
```

For R/RSQLite:

```
years <- seq(from=1987, to=2008)
years_strings <- sapply(years, toString)
fns <- sapply(years_strings, paste, sep="", ".csv.bz2")

# install.packages("RSQLite")
library("RSQLite")
# install.packages("str_pad")
library("stringr")
```

```

my_path <- "~/\"
setwd(my_path)

database_filename = "FlightDatabase.sqlite"
ptm <- proc.time()
for (i in seq(length(fns)))
{
  print(fns[[i]])
  my_bz <- bzfile(fns[[i]])
  my_csv <- read.csv(my_bz,header=TRUE)

  my_csv[is.na(my_csv)] <- 0.1234
  my_csv$DepTime <- substr(str_pad(my_csv$DepTime, 4, pad="0"), 1, 2)

  drv <- dbDriver("SQLite")
  db <- dbConnect(drv, dbname = database_filename)

  dbWriteTable(conn = db, name = "flight_info",
               value = my_csv, row.names = FALSE, append = TRUE)
}
proc.time() - ptm

dbSendQuery(db, "delete from flight_info where DepDelay==0.1234")
dbSendQuery(db, "delete from flight_info where DepTime is '0.'")

# 1a)

file.info(database_filename)

# 2b)
ptm <- proc.time()
x <- fetch(dbSendQuery(db, "select UniqueCarrier, Origin,
                        Dest, Month, DayOfWeek,
                        DepTime,
                        SUM(CASE WHEN DepDelay > 60 THEN 1 ELSE 0 END) as DelayedCounts,
                        Count(*) as TotalFlightCounts,
                        CAST(SUM(CASE WHEN DepDelay > 60 THEN 1.0 ELSE 0.0 END) AS FLOAT) / Count(*)
                        as DelayFraction
                        from flight_info group by
                        UniqueCarrier, Origin, Dest, Month, DayOfWeek, DepTime
                        order by DelayFraction desc"),n=-1)
y <- fetch(dbSendQuery(db, "select UniqueCarrier, Origin,
                        Dest, Month, DayOfWeek,
                        DepTime,
                        SUM(CASE WHEN DepDelay > 90 THEN 1 ELSE 0 END) as DelayedCounts,
                        Count(*) as TotalFlightCounts,
                        CAST(SUM(CASE WHEN DepDelay > 90 THEN 1.0 ELSE 0.0 END) AS FLOAT) / Count(*)
                        as DelayFraction
                        from flight_info group by
                        UniqueCarrier, Origin, Dest, Month, DayOfWeek, DepTime
                        order by DelayFraction desc"),n=-1)
z <- fetch(dbSendQuery(db, "select UniqueCarrier, Origin,

```

```

        Dest, Month, DayOfWeek,
        DepTime,
        SUM(CASE WHEN DepDelay > 180 THEN 1 ELSE 0 END) as DelayedCounts,
        Count(*) as TotalFlightCounts,
        CAST(SUM(CASE WHEN DepDelay > 180 THEN 1.0 ELSE 0.0 END) AS FLOAT) / Count(*)
        as DelayFraction
    from flight_info group by
    UniqueCarrier, Origin, Dest, Month, DayOfWeek, DepTime
    order by DelayFraction desc"),n=-1)

proc.time() - ptm

head(x,n=10)
head(y,n=10)
head(z,n=10)

```

Listing 4: Timing for the loading, initial (non-indexed) query

```

> proc.time() - ptm
      user      system elapsed
3041.288      64.808 3124.991

> proc.time() - ptm
      user      system elapsed
1846.000     141.804 2270.214

> head(x,n=10)
  UniqueCarrier Origin Dest Month DayOfWeek DepTime DelayedCounts
1             9E   ABE  DTW     1         1      08             1
2             9E   ABE  DTW     1         1     14             1
3             9E   ABE  DTW     1         2     18             1
4             9E   ABE  DTW     1         3     14             1
5             9E   ABE  DTW     1         3     18             1
6             9E   ABE  DTW     1         5     17             1
7             9E   ABE  DTW     1         5     23             1
8             9E   ABE  DTW     1         7     14             1
9             9E   ABE  DTW     2         2      07             1
10            9E   ABE  DTW     2         2      08             1
  TotalFlightCounts DelayFraction
1                  1             1
2                  1             1
3                  1             1
4                  1             1
5                  1             1
6                  1             1
7                  1             1
8                  1             1
9                  1             1
10                 1             1

> head(y,n=10)
  UniqueCarrier Origin Dest Month DayOfWeek DepTime DelayedCounts
1             9E   ABE  DTW     1         1      08             1
2             9E   ABE  DTW     1         1     14             1
3             9E   ABE  DTW     1         2     18             1
4             9E   ABE  DTW     1         3     14             1

```

```

5           9E      ABE   DTW      1           3       18           1
6           9E      ABE   DTW      1           5       23           1
7           9E      ABE   DTW      1           7       14           1
8           9E      ABE   DTW      2           2       08           1
9           9E      ABE   DTW      2           2       15           1
10          9E      ABE   DTW      2           7       17           1
  TotalFlightCounts DelayFraction
1                1              1
2                1              1
3                1              1
4                1              1
5                1              1
6                1              1
7                1              1
8                1              1
9                1              1
10               1              1
> head(z,n=10)
  UniqueCarrier Origin Dest Month DayOfWeek DepTime DelayedCounts
1           9E    ABE   DTW     1          3       18             1
2           9E    ABE   DTW     1          5       23             1
3           9E    ABE   DTW     2          2       15             1
4           9E    ABE   DTW     6          1       10             1
5           9E    ABE   DTW     6          2       21             1
6           9E    ABE   DTW     9          7       20             1
7           9E    ABE   DTW    12          1       20             1
8           9E    ABE   DTW    12          3       11             1
9           9E    ABE   DTW    12          5       21             1
10          9E    ABE   DTW    12          7       15             1
  TotalFlightCounts DelayFraction
1                1              1
2                1              1
3                1              1
4                1              1
5                1              1
6                1              1
7                1              1
8                1              1
9                1              1
10               1              1

```

2.3

We perform the same calculation, just using python instead of R. It gives the right answer ONLY if the floating point is cast first to Digit. If the code chunk is run “as is” from the problem statement, it sums to 1.

2.4

We add a index, which speeds up our calculation precipitously! Note, I ran this in the middle of the night, and I think the process got hung in an odd way, but the user and system time were much faster.

```

dbSendQuery(db, "create index delay_index on flight_info
                (UniqueCarrier, Origin, Dest, Month,
                 DayOfWeek, DepTime)")
# dbSendQuery(db, "drop index delay_index")

ptm <- proc.time()
x <- fetch(dbSendQuery(db, "select UniqueCarrier, Origin,
                             Dest, Month, DayOfWeek,
                             DepTime,
                             SUM(CASE WHEN DepDelay > 60 THEN 1 ELSE 0 END) as DelayedCounts,
                             Count(*) as TotalFlightCounts,
                             CAST(SUM(CASE WHEN DepDelay > 60 THEN 1.0 ELSE 0.0 END) AS FLOAT) / Count(*)
                             as DelayFraction
                             from flight_info group by
                             UniqueCarrier, Origin, Dest, Month, DayOfWeek, DepTime
                             order by DelayFraction desc"),n=-1)

y <- fetch(dbSendQuery(db, "select UniqueCarrier, Origin,
                             Dest, Month, DayOfWeek,
                             DepTime,
                             SUM(CASE WHEN DepDelay > 90 THEN 1 ELSE 0 END)
                             as DelayedCounts,
                             Count(*) as TotalFlightCounts,
                             CAST(SUM(CASE WHEN DepDelay > 90 THEN 1.0 ELSE 0.0 END) AS FLOAT) / Count(*) as I
                             from flight_info group by
                             UniqueCarrier, Origin, Dest, Month, DayOfWeek, DepTime
                             order by DelayFraction desc"),n=-1)

z <- fetch(dbSendQuery(db, "select UniqueCarrier, Origin,
                             Dest, Month, DayOfWeek,
                             DepTime,
                             SUM(CASE WHEN DepDelay > 180 THEN 1 ELSE 0 END) as DelayedCounts,
                             Count(*) as TotalFlightCounts,
                             CAST(SUM(CASE WHEN DepDelay > 180 THEN 1.0 ELSE 0.0 END) AS FLOAT) / Count(*)
                             as DelayFraction
                             from flight_info group by
                             UniqueCarrier, Origin, Dest, Month, DayOfWeek, DepTime
                             order by DelayFraction desc"),n=-1)

proc.time() - ptm

```

Listing 5: Timing for all_preprocess.sh

```

> proc.time() - ptm
      user      system    elapsed
 832.760    139.484 12617.140

```

2.5

Using R, we take our object generated by RSQLite and then subset based on flights with at least 150 entries. Then we view the top 10 for each of the 30, 90, 180 minute delays, respectively.

```
# 2e)

xb <- subset(x, TotalFlightCounts > 149)
yb <- subset(y, TotalFlightCounts > 149)
zb <- subset(z, TotalFlightCounts > 149)

head(xb,n=10)
head(yb,n=10)
head(zb,n=10)
```

Listing 6: Timing for all_preprocess.sh

```
> head(xb,n=10)
  UniqueCarrier Origin Dest Month DayOfWeek DepTime DelayedCounts
1638876      WN   HOU  DAL     6           5        18           36
1666191      WN   HOU  DAL     5           4        21           31
1666878      WN   HOU  DAL     2           5        19           26
1666978      WN   HOU  DAL    10           5        18           33
1732659      WN   HOU  DAL     5           4        19           29
1744973      WN   HOU  DAL    10           5        20           28
1745327      WN   HOU  DAL     6           4        19           28
1749951      WN   DAL  HOU     2           5        21           26
1761811      WN   DAL  HOU     4           5        21           25
1761812      UA   LAX  SFO    10           5        12           23
  TotalFlightCounts DelayFraction
1638876          189      0.1904762
1666191          180      0.1722222
1666878          153      0.1699346
1666978          195      0.1692308
1732659          174      0.1666667
1744973          175      0.1600000
1745327          177      0.1581921
1749951          168      0.1547619
1761811          163      0.1533742
1761812          150      0.1533333

> head(yb,n=10)
  UniqueCarrier Origin Dest Month DayOfWeek DepTime DelayedCounts
1191237      WN   DAL  HOU     6           5        20           16
1221851      WN   HOU  DAL     6           4        19           17
1253250      WN   HOU  DAL     7           7        19           14
1253508      WN   HOU  DAL     5           4        21           16
1257823      WN   HOU  DAL     5           4        19           15
1257836      WN   HOU  DAL     2           5        18           14
1258619      WN   HOU  DAL    10           5        20           15
1258684      WN   HOU  DAL     6           5        21           14
1283348      UA   LAX  SFO    11           7        17           13
1287479      WN   DAL  HOU     6           4        21           12
  TotalFlightCounts DelayFraction
1191237          158      0.10126582
1221851          177      0.09604520
1253250          157      0.08917197
1253508          180      0.08888889
1257823          174      0.08620690
1257836          163      0.08588957
```



```

1258619          175      0.08571429
1258684          164      0.08536585
1283348          157      0.08280255
1287479          150      0.08000000
> head(zb,n=10)
      UniqueCarrier Origin Dest Month DayOfWeek DepTime DelayedCounts
378918           WN   HOU  DAL     7         7      19             5
383602           WN   HOU  DAL     4         5      20             5
397917           WN   HOU  DAL     4         2      21             4
399799           WN   HOU  DAL     7         3      20             4
403164           WN   DAL  HOU     5         4      19             4
403202           WN   HOU  DAL    10         5      20             4
413930           WN   DAL  HOU     6         2      21             3
414237           AA   ORD  DFW    12         4      18             3
415160           UA   SFO  LAX    10         7      16             3
415161           UA   SFO  LAX    12         7      16             3
      TotalFlightCounts DelayFraction
378918             157      0.03184713
383602             167      0.02994012
397917             161      0.02484472
399799             166      0.02409639
403164             173      0.02312139
403202             175      0.02285714
413930             150      0.02000000
414237             153      0.01960784
415160             153      0.01960784
415161             153      0.01960784

```

3

```

# 3)
# install.packages("parallel")
library("parallel")

getDelays <- function(x,s) {
  # print(x)
  s <- toString(s)
  # print(s)
  drv <- dbDriver("SQLite")
  db <- dbConnect(drv, dbname = "Big_v4.sqlite")
  query <- sprintf("select UniqueCarrier, Origin,
                      Dest, Month, DayOfWeek,
                      DepTime,
                      SUM(CASE WHEN DepDelay > %i THEN 1 ELSE 0 END) as DelayedCounts,
                      Count(*) as TotalFlightCounts,
                      CAST(SUM(CASE WHEN DepDelay > %i THEN 1.0 ELSE 0.0 END) AS FLOAT) / Count(*)
                      as DelayFraction
                    from flight_info where Origin like '%s%' group by
                    UniqueCarrier, Origin, Dest, Month, DayOfWeek, DepTime
                    order by DelayFraction desc", x, x, s)

  tmp <- fetch(dbSendQuery(db, query),n=-1)
  return(tmp)
}

```

```

}

alphabet = c("A","B","C","D","E","F","G","H","I","J","K","L",
             "M","N","O","P","Q","R","S","T","U","V","W","X",
             "Y","Z")
times <- c(60, 90, 180)
tlc <- expand.grid(times,alphabet)
names(tlc) <- c("x","s")

ptm <- proc.time()
question_3 <- mcmapply(getDelays, tlc$x, tlc$s, mc.cores=4)
proc.time() - ptm

```

Listing 7: Timing for all_preprocess.sh

```

> proc.time() - ptm
      user   system elapsed
245.192   22.508   375.559

```

4

The preprocessing I used here was basically derived from my solution in ps2, just generalized to work with more than one bzip2 file. The code and a wrapper script is shown below below. The file almost takes about 10 minutes, which is m3.xlarge instance. It is probably worth it, for the case of R, because reading in the files can take a very long time.

Listing 8: all_preprocess.sh

```

myyear=$1
# Extract the header so we can find our columns of interest
bzcatt $myyear.csv.bz2 | head -n 1 > $myyear.header.txt
# We will use the file line coordinates as the proxy for index columns
sed -e 's/./\\n/g' $myyear.header.txt > $myyear.header.nsv
# Our desired headers
for i in "UniqueCarrier" "Origin" "Dest" "Month" "DayOfWeek" "DepTime" "DepDelay"
do
    x=`grep -n ^$i$ $myyear.header.nsv | cut -d':' -f 1`
    v="$v $x"
done
echo $v

# Now $v contains our columns of interest, which we just need
# to separate by commas to use with cut. A sed command will
# accomplish this with ease.

bzcatt $myyear.csv.bz2 | \
    cut -d, -f`echo $v` | \
    sed 's/ /,/g' | bzip2 > $myyear.pp.csv.bz2

```

This script is called with:

Listing 9: preprocess.sh

```

date

```

```
for f in `seq 1987 1 2008`  
do  
    echo $f  
    ~/preprocess.sh $f  
done  
date
```

And the result:

Listing 10: Timing for all_preprocess.sh

```
Mon Nov  2 08:25:32 UTC 2015  
Mon Nov  2 08:35:15 UTC 2015
```