

## Problem Statement

Implement the below blog at your end and send the complete documentation.

[https://drive.google.com/file/d/0B\\_Qjau8wv1KoUThzZ24tT1NsZGs/view?usp=sharing](https://drive.google.com/file/d/0B_Qjau8wv1KoUThzZ24tT1NsZGs/view?usp=sharing)

(Airline\_data\_analysis.odt)

The U.S. Department of Transportation's (DOT) Bureau of Transportation Statistics (BTS) tracks the on-time performance of domestic flights operated by large air carriers. Summary information on the number of on-time, delayed, canceled, and diverted flights appears in DOT's monthly Air Travel Consumer Report, published about 30 days after the month's end, as well as in summary tables posted on this website. Summary statistics and raw data are made available to the public at the time the Air Travel Consumer Report is released.

You can download the datasets from the following links:

[Delayed\\_Flights.csv](#)

## DataSet

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
		Year	Month	DayOfMonth	DayOfWeek	DepTime	CRSDepTime	ArrTime	CRSArrTime	UniqueCarrier	FlightNum	TailNum	ActualElapsedTime	CRSElapsedTime	AirTime	ArrDelay	DepDelay	Origin	Dest	Distance	TaxiIn	TaxiOut
1	2008	2008	1	3	4	2003	1955	2211	2225WN		335N712SW		128	150	116	14		81AD	TPA	810	4	8
2		2008	1	3	4	754	735	1002	1000WN		3231N772SW		128	145	113	2		191AD	TPA	810	4	8
3		2008	1	3	4	628	620	804	750WN		448N428W*		96	90	76	14		81ND	BWI	515	3	17
4		2008	1	3	4	1829	1755	1959	1925WN		3920N464W*		90	90	77	34		341ND	BWI	515	3	10
5		2008	1	3	4	1940	1915	2121	2110WN		378N726SW		101	115	87	11		251ND	JAX	688	4	10
6		2008	1	3	4	1937	1830	2037	1940WN		509N763SW		240	250	230	57		671ND	LAS	1591	3	7
7		2008	1	3	4	706	700	916	915WN		100N690SW		130	135	106	1		61ND	MCO	828	5	19
8		2008	1	3	4	1644	1510	1845	1725WN		1333N334SW		121	135	107	80		941ND	MCO	828	6	8
9		2008	1	3	4	1029	1020	1021	1010WN		2272N263W*		52	50	37	11		91ND	MDW	162	6	9
10		2008	1	3	4	1452	1425	1640	1625WN		675N286W*		228	240	213	15		271ND	PHX	1489	5	16
11	2008	2008	1	3	4	754	745	940	955WN		1144N778SW		226	250	205	15		91ND	PHX	1489	5	16
12		2008	1	3	4	1323	1255	1526	1510WN		4N674AA		123	135	110	16		281ND	TPA	838	4	9
13		2008	1	3	4	1416	1325	1512	1435WN		54N643SW		56	70	49	37		511SP	BWI	220	2	5
14		2008	1	3	4	1657	1625	1754	1735WN		623N724SW		57	70	47	19		321SP	BWI	220	5	5
15		2008	1	3	4	1900	1840	1956	1950WN		717N786SW		56	70	49	6		201SP	BWI	220	2	5
16		2008	1	3	4	1039	1030	1133	1140WN		1244N714CB		54	70	47	-7		91SP	BWI	220	2	5
17		2008	1	3	4	1520	1455	1619	1605WN		2553N394SW		59	70	50	14		251SP	BWI	220	2	7
18		2008	1	3	4	1422	1255	1657	1610WN		188N215W*		155	195	143	47		871SP	FLL	1093	6	6
19		2008	1	3	4	1954	1925	2239	2235WN		1754N243W*		165	190	155	4		291SP	FLL	1093	3	7
20		2008	1	3	4	2107	1945	2334	2230WN		362N798SW		147	165	134	64		821SP	MCO	972	6	7
21	2008	2008	1	3	4	1312	1300	1546	1550WN		1397N247W*		154	170	140	-4		121SP	MCO	972	7	7
22		2008	1	3	4	1449	1430	1715	1720WN		3398N707SA		146	170	134	-5		191SP	MCO	972	6	6
23		2008	1	3	4	1634	1555	1859	1845WN		3480N443W*		145	170	134	14		391SP	MCO	972	5	6
24		2008	1	3	4	1812	1650	1927	1815WN		422N779SW		135	145	118	72		821SP	MDW	765	6	11
25		2008	1	3	4	1127	1105	1235	1230WN		1837N704SW		128	145	114	5		221SP	MDW	765	9	5
26		2008	1	3	4	1424	1355	1531	1520WN		2871N708SW		127	145	113	11		291SP	MDW	765	8	6
27		2008	1	3	4	1326	1230	1559	1530WN		1056N459W*		153	180	143	29		561SP	PBI	1052	5	5
28		2008	1	3	4	1749	1725	2019	2030WN		2175N621SW		150	185	138	-11		241SP	PBI	1052	4	8
29		2008	1	3	4	726	720	958	1020WN		3319N206W*		152	180	140	-22		61SP	PBI	1052	4	8
30		2008	1	3	4	646	640	929	955WN		3657N280W*		163	195	151	-26		61SP	RSW	1101	3	9
31	2008	2008	1	3	4	1153	1140	1428	1440WN		2006N241W*		155	180	143	-12		131SP	TPA	1034	4	8
32		2008	1	3	4	1528	1510	1802	1810WN		3858N200W*		154	180	144	-8		181SP	TPA	1034	4	8
33		2008	1	3	4	1450	1435	1806	1745WN		3244N475W*		136	130	121	21		151AN	BWI	888	7	8
34		2008	1	3	4	2245	1730	2354	1850WN		186N792SW		69	80	59	304		315JAN	HOU	359	3	7
35		2008	1	3	4	2025	1940	2135	2100WN		3154N252W*		70	80	60	35		461JAN	HOU	359	3	7
36		2008	1	3	4	1038	945	1314	1225WN		1035N346SW		96	100	81	49		531JAN	MCO	587	8	7
37		2008	1	3	4	1900	1850	2123	2045WN		205N299W*		143	115	97	38		101AN	MDW	666	40	6
38		2008	1	3	4	948	925	959	940WN		3430N487W*		71	75	59	19		231JAN	BHM	365	3	9
39		2008	1	3	4	646	620	725	655WN		1580N243W*		99	95	77	30		261JAN	BNA	484	6	16
40		2008	1	3	4	1110	1040	1136	1110WN		2195N479W*		86	90	72	26		301JAN	BNA	484	5	9
41	2008	2008	1	3	4	2139	2130	2244	2240WN		378N726SW		65	70	54	4		91JAN	FLL	318	3	8
42		2008	1	3	4	1738	1730	1841	1840WN		3948N487W*		63	70	49	1		81JAN	FLL	318	6	8
43		2008	1	3	4	1813	1735	1936	1905WN		54N643SW		143	150	125	31		81JAN	HOU	816	6	12
44		2008	1	3	4	802	750	1001	955WN		2272N263W*		119	125	104	6		21JAN	JND	688	7	8

## Delayed\_Flights.csv Data Description

There are 29 columns in this dataset

Name	Description
1	Year 1987-2008
2	Month 1-12
3	DayOfMonth 1-31
4	DayOfWeek 1 (Monday) - 7 (Sunday)
5	DepTime actual departure time (local, hhmm)
6	CRSDepTime scheduled departure time (local, hhmm)
7	ArrTime actual arrival time (local, hhmm)
8	CRSArrTime scheduled arrival time (local, hhmm)
9	UniqueCarrier unique carrier code
10	FlightNum flight number
11	TailNum plane tail number
12	ActualElapsedTime in minutes
13	CRSElapsedTime in minutes
14	AirTime in minutes
15	ArrDelay arrival delay, in minutes
16	DepDelay departure delay, in minutes
17	Origin origin IATA airport code
18	Dest destination IATA airport code
19	Distance in miles
20	TaxiIn taxi in time, in minutes
21	TaxiOut taxi out time in minutes
22	Cancelled was the flight cancelled?
23	CancellationCode reason for cancellation (A = carrier, B = weather, C = NAS, D = security)
24	Diverted 1 = yes, 0 = no
25	CarrierDelay in minutes
26	WeatherDelay in minutes
27	NASDelay in minutes
28	SecurityDelay in minutes
29	LateAircraftDelay in minutes

## Solution

### Importing Spark Sql packages

```
scala> import org.apache.spark.sql._
import org.apache.spark.sql._

scala> import sqlContext.implicits._
import sqlContext.implicits._

scala> █
```

### Creating base rdd for delayed flights dataset to be used to solve the queries

- Begin by reading the data file as a text file from the local FS using the Spark context object **sc**. [delayedflights](#)
- The 1st line of **delayedflights**, that is the data description, is saved in [header](#).
- **header** is then used in [delayed\\_flights](#) to eliminate data description not to perform analysis on it.

## Code

```
scala> val delayedflights = sc.textFile("file:///home/acadgild/Downloads/DelayedFlights.csv")
delayedflights: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[56] at textFile at <console>:39

scala> val header = delayedflights.first()
header: String = ,Year,Month,DayOfMonth,DayOfWeek,DepTime,CRSDepTime,ArrTime,CRSArrTime,UniqueCarrier,FlightNum,TailNum,ActualElapsedTime,CRSElapsedTime,AirTime,ArrDelay,DepDelay,Origin,Dest,Distance,TaxiIn,TaxiOut,Cancelled,CancellationCode,Diverted,CarrierDelay,WeatherDelay,NASDelay,SecurityDelay,LateAircraftDelay

scala> val delayed_flights = delayedflights.filter(row => row != header)
delayed_flights: org.apache.spark.rdd.RDD[String] = MapPartitionsRDD[57] at filter at <console>:43

scala> █
```

## Output

Showing only 5 records from delayed flights dataset

```
scala> delayed_flights.take(5).foreach(println)
0,2008,1,3,4,2003.0,1955,2211.0,2225,WN,335,N712SW,128.0,150.0,116.0,-14.0,8.0,IAD,TPA,810,4.0,8.0,0,N,,,,,
1,2008,1,3,4,754.0,735,1002.0,1000,WN,3231,N772SW,128.0,145.0,113.0,2.0,19.0,IAD,TPA,810,5.0,10.0,0,N,,,,,
2,2008,1,3,4,628.0,620,804.0,750,WN,448,N428WN,96.0,90.0,76.0,14.0,8.0,IND,BWI,515,3.0,17.0,0,N,,,,,
4,2008,1,3,4,1829.0,1755,1959.0,1925,WN,3920,N464WN,90.0,90.0,77.0,34.0,34.0,IND,BWI,515,3.0,10.0,0,N,0,2.0,0.0,0.0,0.0,32.0
5,2008,1,3,4,1940.0,1915,2121.0,2110,WN,378,N726SW,101.0,115.0,87.0,11.0,25.0,IND,JAX,688,4.0,10.0,0,N,,,,,

scala> █
```

### Find out the top 5 most visited destinations

- The data is first split by comma, as it is from a comma separated value (csv) file and mapped to only the 18th column (**Dest**) and 1, i.e. (Dest,1)
- The resultant data is then filtered such that there is no blank values in the Destination field
- Next, the resultant data is reduced by the key element (Dest), such that all the data corresponding to each key element are grouped together and added (i.e. all the 1's)
- Next, the resultant data is mapped to get the key with the summed value i.e. (Dest, sum(Dest 1's)). This gives us the count of how many times the Destination has appeared in the data.
- Next, the resultant data is sorted by key **false**, i.e. in descending order.
- Finally, only select the top 5 values from the resultant data.

## Code

```
scala> val most_visited_destn = delayed_flights.map(x => x.split(",")).
  | map(x => (x(18),1)).
  | filter(x => x._1 != null).
  | reduceByKey(_+_).
  | map(x => (x._2, x._1)).
  | sortByKey(false).map(x => (x._2,x._1)).
  | take(5)
most_visited_destn: Array[(String, Int)] = Array((ORD,108984), (ATL,106898), (DFW,70657), (DEN,63003), (LAX,59969))
scala> █
```

---

## Output

```
scala> most_visited_destn.foreach(println)
(ORD,108984)
(ATL,106898)
(DFW,70657)
(DEN,63003)
(LAX,59969)

scala> █
```

---

### Which month has seen the most number of cancellations due to bad weather?

- The data is first split by comma, as it is from a comma separated value (csv) file and filtered to data where,
- 22nd column (**Cancelled**) == 1, signifying that a cancellation had happened *and*
- 23rd column (**CancellationCode**) == B, signifying cancellation due to bad weather
- The data is then mapped to the data that satisfies the above criteria's and 1, i.e. (entire\_row,1)
- Next, the resultant data is reduced by the key element (the entire row-x(2)), such that all the data corresponding to each key element are grouped together and added (i.e. all the 1's)
- Next, the resultant data is mapped to get the **Month(x.\_2)** with the summed value i.e. (Month, sum(Cancels in Month)). This gives us the count of how many times cancellation has happened in that month.
- Next, the resultant data is sorted by key **false**, i.e. in descending order and only the 1st row is selected.

## Code

```
scala> val most_cancellation_month = delayed_flights.map(x => x.split(",")).
  | filter(x => ((x(22).equals("1"))&&(x(23).equals("B")))).
  | map(x => (x(2),1)).
  | reduceByKey(_+_).
  | map(x => (x._2,x._1)).
  | sortByKey(false).
  | map(x => (x._2,x._1)).
  | take(1)
most_cancellation_month: Array[(String, Int)] = Array((12,250))
scala> █
```

---

## Output

```
scala> most_cancellation_month.foreach(println)
(12,250)

scala> █
```

---

## Top ten origins with the highest AVG departure delay

- The data is first split by comma, as it is from a comma separated value (csv) file and mapped to only the 17th column (**Origin**) and 16th column (**DepDelay**) respectively
- Both columns are then mapped to (column, 1). i.e. (Origin,1) and (DepDelay,1)
- Next, the resultant data is reduced by the key element for both Origin and DepDelay, such that all the data corresponding to each key element are grouped together and added (i.e. all the 1's for each of them)
- Next, take the Origin as the **sum** and DepDelay as the **count** and get the **avg** for each row and map the results as (Origin, avg)
- Finally, the resultant data is sorted by key **false**, i.e. in Desc order and 1st 10 rows are selected.

## Code

```
scala> val avg = delayed_flights.map(x => x.split(",")).
  | map(x => (x(17),x(16).toDouble)).
  | mapValues(_._1).
  | reduceByKey((x, y) => (x._1 + y._1, x._2 + y._2)).
  | mapValues{ case (sum, count) => (1.0 * sum)/count}.
  | map(x => (x._2,x._1)).
  | sortByKey(false).
  | map(x => (x._2,x._1)).
  | take(10)
avg: Array[(String, Double)] = Array((CMX,116.1470588235294), (PLN,93.76190476190476), (SPI,83.84873949579831), (ALO,82.22580
64516129), (MQT,79.55665024630542), (ACY,79.3103448275862), (MOT,78.66165413533835), (HHH,76.53005464480874), (EGE,74.1289198
6062718), (BGM,73.15533980582525))
scala>
```

## Output

```
scala> avg.foreach(println)
(CMX,116.1470588235294)
(PLN,93.76190476190476)
(SPI,83.84873949579831)
(ALO,82.2258064516129)
(MQT,79.55665024630542)
(ACY,79.3103448275862)
(MOT,78.66165413533835)
(HHH,76.53005464480874)
(EGE,74.12891986062718)
(BGM,73.15533980582525)

scala> █
```

## Which route (origin & destination) has seen the maximum diversion?

- The data is first split by comma, as it is from a comma separated value (csv) file and filtered to data where, the 24th column (**Diverted**) == 1, signifying that the flight had been diverted
- The resultant data is then mapped such that the 1st column is a concatenation of the column 17 (**Origin**) and column 18 (**Dest**) and 1, i.e. ("Origin,Dest",1)
- Next, the resultant data is reduced by the key element ("Origin,Dest"), such that all the data corresponding to each key element are grouped together and added (i.e. all the 1's)
- The data is then mapped to ("Origin,Dest",1)
- Finally, the resultant data is sorted by key **false**, i.e. in Desc order and 1st 10 rows are selected, with the 1st row as the route with the most diversion.

## Code

```
scala> val diversion = delayed flights.map(x => x.split(",")).
  | filter(x => ((x(24).equals("1")))).
  | map(x => ((x(17)+", "+x(18)),1)).
  | reduceByKey(_+_).
  | map(x => (x._2,x._1)).
  | sortByKey(false).
  | map(x => (x._2,x._1)).
  | take(10)
diversion: Array[(String, Int)] = Array((ORD,LGA,39), (DAL,HOU,35), (DFW,LGA,33), (ATL,LGA,32), (SLC,SUN,31), (ORD,SNA,31), (
MIA,LGA,31), (BUR,JFK,29), (HRL,HOU,28), (BUR,DFW,25))
scala>
```

## Output

```
scala> diversion.foreach(println)
(ORD,LGA,39)
(DAL,HOU,35)
(DFW,LGA,33)
(ATL,LGA,32)
(SLC,SUN,31)
(ORD,SNA,31)
(MIA,LGA,31)
(BUR,JFK,29)
(HRL,HOU,28)
(BUR,DFW,25)

scala>

scala> █
```