## Case Study: Analysis of Airline Delay using Spark

Airline delay is a critical issue affecting both airlines and passengers. In this assignment, you will use

PySpark to analyze a dataset containing information about airline flights and predict flight delays.

Technology: Spark, SQL

Dataset: Flights\_Delay.csv

#### **Questions:**

#### a) Create a new Spark Session with new SparkConfig

from pyspark import SparkConf, SparkContext

from pyspark.sql import SparkSession, HiveContext

config = SparkConf().setAppName("AirlineDelayAnalysis").setMaster('local[4]')
sc = SparkContext.getOrCreate(conf=config)

sc

```
Out[6]: SparkContext

Spark UI
Version
v2.4.8
Master
local[4]
AppName
AirlineDelayAnalysis
```

# b) Create new instance of Spark SQL session and define new DataFrame using Flights\_Delay.csv dataset.

```
flights_df =
spark.read.csv("file:///home/hadoop/Downloads/Flights_Delay.csv",inferSchema =
True, header = True)
flights_df.printSchema()
                  |-- ID: integer (nullable = true)
|-- YEAR: integer (nullable = true)
|-- YEAR: integer (nullable = true)
|-- MONTH: integer (nullable = true)
|-- DAY: integer (nullable = true)
|-- DAY: integer (nullable = true)
|-- AIRLINE: string (nullable = true)
|-- FLIGHT NUMBER: integer (nullable = true)
|-- TAIL_NÜMBER: string (nullable = true)
|-- ORIGIN AIRPORT: string (nullable = true)
|-- DESTINATION AIRPORT: string (nullable = true)
|-- DESTINATION AIRPORT: string (nullable = true)
|-- DEPARTURE_TIME: integer (nullable = true)
|-- DEPARTURE_DELAY: integer (nullable = true)
|-- TAXI_OUT: integer (nullable = true)
|-- WHEELS_OFF: integer (nullable = true)
|-- SCHEDULED_TIME: integer (nullable = true)
|-- SCHEDULED_TIME: integer (nullable = true)
|-- SCHEDULED_TIME: integer (nullable = true)
                   -- AIR TIME: integer (nullable = true)
-- DISTANCE: integer (nullable = true)
-- WHEELS ON: integer (nullable = true)
-- TAXI_IN: integer (nullable = true)
-- SCHEDULED ARRIVAL: integer (nullable = true)
                 -- SCHEDULED ARRIVAL: integer (nullable = true)
-- ARRIVAL TIME: integer (nullable = true)
-- ARRIVAL DELAY: integer (nullable = true)
-- DIVERTED: integer (nullable = true)
-- CANCELLED: integer (nullable = true)
-- CANCELLATION REASON: string (nullable = true)
-- AIR SYSTEM DELAY: integer (nullable = true)
-- SECURITY DELAY: integer (nullable = true)
-- AIRLINE DELAY: integer (nullable = true)
-- AIRLINE DELAY: integer (nullable = true)
-- LATE AIRCRAFT DELAY: integer (nullable = true)
-- WEATHER DELAY: integer (nullable = true)
c) Create table Spark HIVE table flights_table
spark.sql("create database if not exists flights_db").show()
spark.sql("show databases").show()
spark.sql("use flights_db").show()
          |databaseName|
                   banking db|
                            default
                   flights db
spark.sql("""
CREATE TABLE IF NOT EXISTS flights_table (
  ID INT,
  YEAR INT,
  MONTH INT,
   DAY INT,
```

DAY\_OF\_WEEK INT,

```
AIRLINE STRING,
FLIGHT_NUMBER STRING,
TAIL_NUMBER STRING,
ORIGIN_AIRPORT STRING,
DESTINATION_AIRPORT STRING,
SCHEDULED_DEPARTURE INT,
DEPARTURE_TIME INT,
DEPARTURE_DELAY INT,
TAXI_OUT INT,
WHEELS_OFF INT,
SCHEDULED_TIME INT,
ELAPSED_TIME INT,
AIR_TIME INT,
DISTANCE INT,
WHEELS_ON INT,
TAXI_IN INT,
SCHEDULED_ARRIVAL INT,
ARRIVAL_TIME INT,
ARRIVAL_DELAY INT,
DIVERTED INT,
CANCELLED INT,
CANCELLATION_REASON STRING,
AIR_SYSTEM_DELAY INT,
SECURITY_DELAY INT,
AIRLINE_DELAY INT,
LATE_AIRCRAFT_DELAY INT,
WEATHER DELAY INT
```

```
ROW FORMAT DELIMITED
FIELDS TERMINATED BY "
STORED AS TEXTFILE
TBLPROPERTIES ("skip.header.line.count" = "1")
df_selected = flights_df.select('ID', 'YEAR', 'MONTH', 'DAY', 'DAY_OF_WEEK', 'AIRLINE',\
              'FLIGHT_NUMBER', 'TAIL_NUMBER', 'ORIGIN_AIRPORT',
'DESTINATION AIRPORT',\
              'SCHEDULED_DEPARTURE', 'DEPARTURE_TIME', 'DEPARTURE_DELAY',
'TAXI OUT',\
              'WHEELS_OFF', 'SCHEDULED_TIME', 'ELAPSED_TIME', 'AIR_TIME',
'DISTANCE',\
              'WHEELS_ON', 'TAXI_IN', 'SCHEDULED_ARRIVAL', 'ARRIVAL_TIME',
'ARRIVAL_DELAY',\
              'DIVERTED', 'CANCELLED', 'CANCELLATION REASON',
'AIR_SYSTEM_DELAY',\
              'SECURITY_DELAY', 'AIRLINE_DELAY', 'LATE_AIRCRAFT_DELAY',
'WEATHER_DELAY')
df_selected.createOrReplaceTempView("flights")
spark.sql("""
insert into table flights_table
select * from flights
""")
d) Describe the table schema & Dataset show top 10 rows of Dataset
spark.sql("describe formatted flights_table").show(10)
```

```
col_name|data_type|comment|
                  IDI
                           intl
                                  nulli
                YEAR |
                           int
                                  null
               MONTH
                           intl
                                  nulli
                 DAY
                           int
                                  null
         DAY OF WEEK
                           int
                                  null
                                  null
             AIRLINE
                        string
       FLIGHT NUMBER
                        string
                                  null
         TAIL NUMBER
                        string
                                  null
      ORIGIN AIRPORT
                        string
                                  null
DESTINATION AIRPORT
                        string|
                                  null
only showing top 10 rows
```

# e) Apply Query performance optimization techniques like – creating Partitioning DataFrame by

# caching

flights\_cached = flights\_df.cache()

flights\_cached

Out[19]: DataFrame[ID: int, YEAR: int, MONTH: int, DAY: int, DAY OF WEEK: int, AIRLINE: string, FLIGHT NUMBER: int, TAIL\_NUMB
ER: string, ORIGIN\_AIRPORT: string, DESTINATION\_AIRPORT: string, SCHEDULED\_DEPARTURE: int, DEPARTURE\_TIME: int, DEPA
RTURE DELAY: int, TAXI OUT: int, WHEELS OFF: int, SCHEDULED TIME: int, ELAPSED TIME: int, AIR TIME: int, DISTANCE: i
 nt, WHEELS\_ON: int, TAXI IN: int, SCHEDULED ARRIVAL: Int, ARRIVAL TIME: int, ARRIVAL DELAY: int, DIVERTED: int, CANC
ELLED: int, CANCELLATION\_REASON: string, AIR\_SYSTEM\_DELAY: int, SECURITY\_DELAY: int, AIRLINE\_DELAY: int, LATE\_AIRCRA
FT\_DELAY: int, WEATHER\_DELAY: int]

# parquet

flights df.write.parquet("file:///home/hadoop/Downloads/Flights/")

parquet\_flight = spark.read.parquet("file:///home/hadoop/Downloads/Flights/part-00000-9b7b224a-b75c-4477-a13c-90a2734dc18b-c000.snappy.parquet")

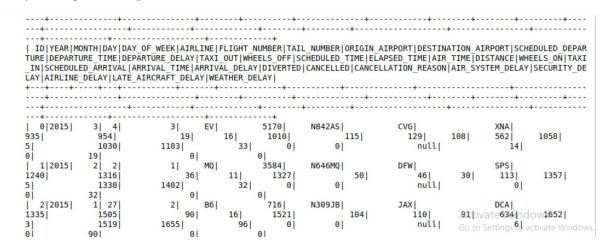
parquet\_flight.show()

```
| ID|YEAR|MONTH|DAY|DAY_OF_WEEK|AIRLINE|FLIGHT_NUMBER|TAIL_NUMBER|ORIGIN_AIRPORT|DESTINATION_AIRPORT|SCHEDULED_DEPAR
TURE|DEPARTURE_TIME|DEPARTURE_DELAY|TAXI_OUT|WHEELS_OFF|SCHEDULED_TIME|ELAPSED_TIME|AIR_TIME|DISTANCE|WHEELS_ON|TAXI
_IN|SCHEDULED_ARRIVAL|ARRIVAL_TIME|ARRIVAL_DELAY|DIVERTED|CANCELLED|CANCELLATION_REASON|AIR_SYSTEM_DELAY|SECURITY_DE
LAY AIRLINE_DELAY | LATE_AIRCRAFT_DELAY | WEATHER_DELAY |
0 | 2015 |
                                                                                  N842AS
                    9541
                                      1103|
                                                                                                              1291
935
                                                        16|
                                                                                            115|
                                                                                                                                                    1058
                                                                                         0 |
                                                                                                                null
                    1030|
                                                             33|
                 19|
2| 2|
                                        1|
36|
                                                  MQ| 11|
   1 | 2015 |
                                                                    35841
                                                                                  N646MQ
                                                                                                                                       SPSI
                                                                                                                 46|
                                                                                                                             30|
                    1316
                                                                                                                                       113|
                                                             32|
                                       1402|
                    1330
                                                                                                                 null
                                                                                                                            Activate Woldows
                                                                  0|
716|
                                              0 |
                  321
                                                                                                                             Go to Settings to activate Windows
   2|2015|
                                         21
                                                   B6 I
                                                                                  N309JBI
                                                                                                           JAXI
```

# partitioning

partitioned\_flights = flights\_df.repartition(3)

partitioned\_flights.write.parquet("file:///home/hadoop/Downloads/Flights/Repartition")
parquet\_flight = spark.read.parquet("file:///home/hadoop/Downloads/Flights/part00000-9b7b224a-b75c-4477-a13c-90a2734dc18b-c000.snappy.parquet")
parquet\_flight.show()



#### # pushdown

parquet\_flight.printSchema()

parquet\_flight.filter(parquet\_flight["MONTH"]!=1).show()

| ID|YEAR|MONTH|DAY|DAY OF WEEK|AIRLINE|FLIGHT NUMBER|TAIL NUMBER|ORIGIN AIRPORT|DESTINATION AIRPORT|SCHEDULED DEPAR TURE DEPARTURE TIME DEPARTURE DELAY TAXÍ OUT WHEELS OFF SCHEDULED TIME ELAPSED TIME AIR TIME DISTANCE WHEELS ON TAXI IN|SCHEDULED\_ARRIVAL|ARRIVAL\_TIME|ARRIVAL\_DÉLAY|DIVERTÉD|CANCELLED|CANCELLATION\_REASON|AIR\_SYSTEM\_DELAY|SECÜRITY\_DE LAY AIRLINE DELAY LATE AIRCRAFT DELAY WEATHER DELAY 3 4 0 | 2015 | 3| 5170 N842AS| CVGI XNA I 954 9351 191 16 115 1291 1081 562 10581 1010 0| null 5| 1030 1103 33 0 19| 0 | 1 | 2015 | 2 MQ | 35841 N646MQ | SPSI 1240 1316 461 13571 361 111 1327 501 113 5| 1330 0 | nuil 0 0 | 51 4 | 2015 | 2 EVI 25| 55841 N851AS Astivate 164 indow 1349 | Go to Setting 119 activate Windows 1255 1250 -51 481 621 1343 1352 0 | null 3|

#### f) Average arrival delay caused by airlines

spark.sql("""

select airline, avg(airline\_delay) Average\_airline\_delay from flights\_table group by airline """).show()

```
|airline|Average_airline_delay|
+----+
              13.272425249169435
       NKI
              21.797598627787306
              21.38169014084507
       B6 i
               19.62044817927171
             24.849923430321592
17.7892749244713
       DL
       F9|
US|
MQ|
HA|
AS|
VX|
              14.47266881028939
18.03059975520196
              15.768812330009066
              24.591836734693878
              13.192118226600986
               10.85576923076923
             14.513646532438479
```

## g) Days of months with respected to average of arrival delays

```
spark.sql("""
select Month,Day,avg(arrival_delay) Average_of_arrival_delay from flights_table
group by month,day
order by month,day
""").show()
```

	+	
1	1	5.93913043478260
1	2	9.6286031042128
1	3	26.09056122448979
1	4	31.8723926380368
1	5	19.85880980163360
1	6	21.01891551071878
1	7	13.59193954659949
1	8	13.30111524163568
1	9	11.0687285223367
1	10	3.107352941176470
1	11	9.54498044328552
1	12	19.9133165829145
1	13	5.1511216056670
1	14	1.024154589371980
1	15	-0.509988249118683
1	16	-0.583333333333333
1	17	-3.609
1	18	-0.756035578144853
1	19	-2.30392156862745

#### h) Arrange weekdays with respect to the average arrival delays caused

```
spark.sql("""
select Day_of_Week,avg(arrival_delay) Average_arrival_delay
from flights_table
group by day_of_week
order by avg(arrival_delay) desc
""").show()
```

#### i) Arrange Days of month as per cancellations done in Descending

```
spark.sql("""
select Day,Month,count(Cancelled) Cancellations from flights_table
where cancelled = 1
group by Day,Month
order by count(Cancelled) desc
""").show()
```

Day M	onth Cance	ellations
++-	+	
5	3	348
2	2	308
27	1	298
1	2	244
1 1	3	180
15	2	158
26	1	152
4	3	150
23	2	134
9	2	120
16	2	120
17	2	116
21	2	108
25	2	104
28	2	102
28	11	94

#### j) Find Top 10 busiest airports with respect to day of week

SELECT destination\_airport AS airport,

```
spark.sql("""

SELECT rank() OVER (PARTITION BY t.day_of_week ORDER BY sum(t.flight_count) DESC)

AS rank,

t.airport,

t.day_of_week,

sum(t.flight_count) AS total_flights

FROM (
```

```
day_of_week,
    count(*) AS flight_count

FROM flights_table

GROUP BY destination_airport, day_of_week

UNION ALL

SELECT origin_airport AS airport,
    day_of_week,
    count(*) AS flight_count

FROM flights_table

GROUP BY origin_airport, day_of_week

) t

GROUP BY t.airport, t.day_of_week

ORDER BY t.day_of_week, total_flights DESC

LIMIT 10

""").show()
```

K a.	cipor ciday_c	of_week tota	
1	ATLI	1	2212
2	ORD	1	1688
3 j	DFW	1	1636
4	LAX	1	1262
5 j	DEN	1	1226
5 j	IAH	1	988
7 j	PHX	1	970
ві	SF0	1	932
9 j	LAS	1	796
j	MSP	1	764

#### k) Finding airlines that make the maximum number of cancellations

```
spark.sql("""
select Airline,count(cancelled) No_Of_Cancellations from flights_table
where cancelled = 1
group by airline
```

order by count(cancelled) desc

""").show()

irline No Of Ca	ncellations
ri cine ino oi ca	incertations
MQ	828
WN	716
EV	624
AAI	482
DL	354
US	338
001	306
B6	290
UAI	244
NK	42
VX	26
ASI	24
F9	22
HA	6

## l) Find and order airlines in descending that make the most number of diversions

spark.sql("""

select airline, count (diverted) as total\_diversion

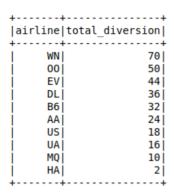
from flights\_table

where diverted = 1

group by airline

order by total\_diversion desc

""").show()



## m) Finding days of month that see the most number of diversion

spark.sql("""

select month,day,count(diverted) as total\_diversion

from flights\_table

where diverted = 1
group by month,day
order by total\_diversion desc

""").show()

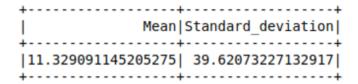
++	+	+
month	day	total_diversion
1 21	2	18
3	5	14
2	14	14
3	1	14
3	2	12
3	4	12
1	18	10
2	23	10
1	11	10
1	7	10
1	30	10
2	1	10
2	21	8
1 1	8	8
2	9	8
1 1	4	6
2	4	6
1 21	281	61

## n) Calculating mean and standard deviation of departure delay for all flights in minutes

spark.sql("""

select mean(departure\_delay) as Mean, std(departure\_delay) as Standard\_deviation from flights\_table

""").show()



## o) Calculating mean and standard deviation of arrival delay for all flights in minutes

spark.sql("""

select mean(arrival\_delay) as Mean, std(arrival\_delay) as Standard\_deviation from flights\_table

""").show()

## p) Finding all diverted Route from a source to destination Airport & Diverted route is the most diverted

```
spark.sql("""
select origin_airport, destination_airport, count(diverted) as No_of_Diversion
from flights_table
where diverted = 1
group by origin_airport, destination_airport
order by No_of_Diversion
""").show()
```

origin_airport des	tination_airport	No_of_Diversion
SF0	BOI	2
į ATLį	GSP	2
I SNA	SF0	2
j FLL j	BWI	2
BOS	LAS	2
SBP	SF0	2
SLC	RDM	2
j MCO j	PVD	2
j FLL	PVD	2
j SLC	SUN	2
CLT	MIA	2
cos	ORD	2
KOA	SF0	2
[ EWR	STL	2
MCO	BWI	2
ATL	ASE	2
į ATLį	GTR	2
[ CAK	LGA	2
		- 1

q) Finding AIRLINES with its total flight count, total number of flights arrival delayed by more than 30 Minutes, % of such flights delayed by more than 30 minutes when it is not Weekends with minimum count of flights from Airlines by more than 10. Also Exclude some of Airlines AK, HI, PR, VI and arrange output in descending order by % of such count of flights.

```
spark.sql("""
WITH airline_query AS (
```

```
SELECT
   f1.airline,
   COUNT(f1.flight_number) AS flights_travelled,
   COUNT(DISTINCT f1.flight_number) AS flights_each_airline,
   SUM(CASE WHEN f1.arrival_delay > 30 THEN 1 ELSE 0 END) AS
total_flights_delayed_by_30mins
  FROM flights_table f1
 GROUP BY f1.airline
)
SELECT
 f.airline,
  q.flights_travelled,
  q.flights_each_airline,
  q.total_flights_delayed_by_30mins,
  SUM(
   CASE
     WHEN f.arrival_delay > 30
       AND f.day_of_week NOT IN (6, 7)
       AND q.flights_each_airline > 10
       AND f.airline NOT IN ('AK', 'HI', 'PR', 'VI')
     THEN 1
     ELSE 0
   END
 ) / q.flights_travelled * 100 AS per_of_fl_delay_30
FROM flights_table f
INNER JOIN airline_query q
  ON f.airline = q.airline
GROUP BY
```

```
f.airline,
  q.flights_travelled,
  q.flights_each_airline,
  q.total_flights_delayed_by_30mins
ORDER BY per_of_fl_delay_30 DESC
""").show()
  |airline|flights travelled|flights each_airline|total_flights_delayed_by_30mins|per_of_fl_delay_30|
        F9I
                        15881
                                              2941
                                                                             396 | 17.506297229219143 |
        MOI
                        7004
                                              7401
                                                                            1550 | 17.16162193032553
        B61
                        5096
                                              770
                                                                             970 | 14.128728414442701 |
                        2096
                                              261
                                                                             372 13.263358778625955
                       11832
                                             1685
                                                                             1748 11.240703177822853
                                             1462
                                                                             1718 | 11.089698668535389
        00
                       11416
                                                                            1306 | 10.57221867687726
        UA
                        94021
                                             1167
        AA
                       10500
                                             1129
                                                                            1400 | 9.219047619047618 |
        VX
                        1146
                                              163
                                                                             134 8.202443280977311
        US
                       7850
                                             836
                                                                             904 7.898089171974522
                                                                             1492 7.410189009888597
                       15978
                                             2022
        WN
                       23476
                                             2847
                                                                             2470
                                                                                  7.40330550349293
                                                                                   4.03530895334174
                        3172
                                              435
                                                                             2001
        AS
```

r) Finding AIRLINES with its total flight count with total number of flights departure delayed by less than 30 Minutes, % of such flights delayed by less than 30 minutes when it is Weekends with minimum count of flights from Airlines by more than 10. Also Exclude some of Airlines AK, HI, PR, VI and arrange output in descending order by % of such count of flights.

76 3 . 1855 955 678 670 362 |

210

```
spark.sql("""
WITH airline_query AS (
SELECT
f1.airline,
COUNT(f1.flight_number) AS flights_travelled,
COUNT(DISTINCT f1.flight_number) AS flights_each_airline,
SUM(CASE WHEN f1.departure_delay < 30 THEN 1 ELSE 0 END) AS total_flights_delayed_less_30_mins
FROM flights_table f1
GROUP BY f1.airline
)
SELECT
```

HAI

1444

```
f.airline,
  q.flights_travelled,
  q.flights_each_airline,
  q.total_flights_delayed_less_30_mins,
  SUM(
   CASE
     WHEN f.departure_delay < 30
        AND f.day_of_week IN (6, 7)
        AND q.flights_each_airline > 10
        AND q.airline NOT IN ('AK', 'HI', 'PR', 'VI')
     THEN 1
     ELSE 0
   END
 ) / q.flights_travelled * 100 AS per_of_fl_delay_30
FROM flights_table f
INNER JOIN airline_query q
  ON f.airline = q.airline
GROUP BY
 f.airline,
  q.flights_travelled,
  q.flights_each_airline,
  q.total_flights_delayed_less_30_mins
ORDER BY per_of_fl_delay_30 DESC
""").show()
```

airline flight	ts_travelled flights	_each_airline tota	l_flights_delayed_less_30_mins per_of_fl_delay_30
+			
AS	3172	435	2936   25.97730138713745
HA	1444	210	1384   24.792243767313018
NK	2096	261	1678   24.141221374045802
AA	10500	1129	8684   23.123809523809523
j DL j	15978	2022	14020 22.70622105394918
į VX į	1146	163	980 22.5130890052356
į WN į	23476	2847	19890 22.456977338558527
j usj	7850	836	6712 22.089171974522294
j 00 j	11416	1462	9472 21.79397337070778
j B6 j	5096	770	3894 21.31083202511774
į EVį	11832	1685	9638 20.33468559837728
i UAİ	9402	1167	7806 20.2084662837694
i MQ i	7004	740	4886 17.76127926898915
F9	1588	294	1170 16.750629722921914
<b>‡</b>			Activate W

# s) When is the best time of day/day of week/time of a year to fly with minimum delays?

spark.sql("""

**SELECT** 

day\_of\_week,

AVG(COALESCE(arrival\_delay, 0) + COALESCE(departure\_delay, 0)) AS total\_delay

FROM flights\_table

GROUP BY day\_of\_week

ORDER BY total\_delay

""").show()

total_delay	day_of_week
	+
14.154561301568855	6
14.233472149921916	3
15.901192887688499	5
17.397617629541394	4
18.47514450867052	2
22.593458540948	7
23.541281180466097	1

t) Which airlines are best airline to travel considering number of cancellations, arrival, departure delays and all reasons affecting performance of airline industry. spark.sql("""

**SELECT** 

Airline,

COUNT(flight\_number) AS total\_flights,

SUM(CASE WHEN cancelled = 1 THEN 1 ELSE 0 END) AS total\_cancellations,

AVG(COALESCE(arrival\_delay, 0)) AS avg\_arrival\_delay,

AVG(COALESCE(departure\_delay, 0)) AS avg\_departure\_delay,

SUM(COALESCE(air\_system\_delay, 0) +

COALESCE(security\_delay, 0) +

COALESCE(late\_aircraft\_delay, 0) +

COALESCE(airline\_delay, 0) +

COALESCE(weather\_delay, 0)

) AS total\_delay\_reasons

FROM flights\_table

**GROUP BY airline** 

""").show()

+					
Airline total	_flights tota	al_cancellations	avg_arrival_delay	avg_departure_delay	total_delay_reasons
+		+			······
UA	9402	244	6.5120187194214	13.940438204637312	120396
NK	2096	42	13.92175572519084	15.283396946564885	37624
AA	10500	482	7.9824761904761905	10.987428571428572	141970
EV	11832	624	10.269776876267748	10.95368492224476	173040
B6	5096	290	13.076530612244898	15.154238618524333	92368
DL	15978	354	2.7457754412316935	9.71861309300288	168084
00	11416	306	9.83812193412754	11.302733006306937	166506
F9	1588	22	23.769521410579344	23.188916876574307	42370
US	7850	338	5.70624203821656	7.484076433121019	84370
MQ	7004	828	16.93061107938321	15.098515134209023	138706
HA	1444	6	4.049861495844875	1.1786703601108033	8144
AS	3172	24	-1.5201765447667086	2.320302648171501	20976
VX	1146	26	5.012216404886562	9.633507853403142	13422
WN	23476	716	3.574033055034929	9.810103935934572	Activate Wind

Go to Settings to act