

Stage 3 - Database Design

Part 1:

Data Definition Language (DDL) commands:

Airports table

```
CREATE TABLE airports (
    IATA_CODE VARCHAR(3) PRIMARY KEY,
    AIRPORT VARCHAR(255),
    CITY VARCHAR(255),
    STATE VARCHAR(2)
);
```

Airlines table

```
CREATE TABLE airlines (
    IATA_CODE VARCHAR(2) PRIMARY KEY,
    AIRLINE VARCHAR(255)
);
```

AirlineRating table

```
Create table AirlineRating(
    AIRLINE VARCHAR(255) PRIMARY KEY,
    ReliabilityRating INT,
    OnTimePerformanceRating INT,
    FOREIGN KEY(AIRLINE) REFERENCES airlines(IATA_CODE)
);
```

Flight table

```
CREATE TABLE flights(
    DATE VARCHAR(255),
    DAY_OF_WEEK INT,
    AIRLINE VARCHAR(3),
    FLIGHT_NUMBER INT,
    ORIGIN_AIRPORT VARCHAR(3),
```

```
    DESTINATION_AIRPORT VARCHAR(3),  
    PRIMARY KEY(DATE, AIRLINE, FLIGHT_NUMBER),  
    FOREIGN KEY(ORIGIN_AIRPORT) REFERENCES airports(IATA_CODE),  
    FOREIGN KEY(DESTINATION_AIRPORT) REFERENCES  
airports(IATA_CODE),  
    FOREIGN KEY(AIRLINE) REFERENCES airlines(IATA_CODE)  
);
```

Delay table

```
Create Table delays(  
    DATE VARCHAR(255),  
    AIRLINE VARCHAR(3) ,  
    FLIGHT_NUMBER INT ,  
    SCHEDULED_DEPARTURE INT,  
    DEPARTURE_TIME INT,  
    DEPARTURE_DELAY INT,  
    SCHEDULED_TIME INT,  
    ELAPSED_TIME INT,  
    AIR_TIME INT,  
    DISTANCE INT,  
    SCHEDULED_ARRIVAL INT,  
    ARRIVAL_TIME INT,  
    ARRIVAL_DELAY INT,  
    DIVERTED INT,  
    CANCELLED INT,  
    CANCELLATION_REASON VARCHAR(1),  
    AIR_SYSTEM_DELAY INT,  
    SECURITY_DELAY INT,  
    AIRLINE_DELAY INT,  
    LATE_AIRCRAFT_DELAY INT,  
    WEATHER_DELAY INT,  
    PRIMARY KEY(DATE, AIRLINE, FLIGHT_NUMBER),  
    FOREIGN KEY(DATE, AIRLINE,FLIGHT_NUMBER) REFERENCES  
flights(DATE,      AIRLINE,FLIGHT_NUMBER)  
);
```

Inserting Data & Table counts

In GCP, we uploaded the CSV data into buckets, and imported the values into SQL directly. The original flights dataset only had ~300 rows for the airports table and 14 rows for the airlines table. We used a Python program to add dummy data to pad our table with the required 1000 rows. We also removed duplicate entries from the delays table, which brought down the count from ~1M to 750K.

The final count of rows in each table is as follows:

1. airports: 1001 rows
2. Flights: 749794 rows
3. delays: 749794 rows
4. airlines: 1000 rows
5. AirlineRating: 1000 rows

```
mysql> SELECT COUNT(*) FROM airports;
+-----+
| COUNT(*) |
+-----+
|      1001 |
+-----+
1 row in set (0.01 sec)

mysql> SELECT COUNT(*) FROM flights;
+-----+
| COUNT(*) |
+-----+
|     749794 |
+-----+
1 row in set (0.80 sec)

mysql> SELECT COUNT(*) FROM delays;
+-----+
| COUNT(*) |
+-----+
|    749794 |
+-----+
1 row in set (1.29 sec)

mysql> SELECT COUNT(*) FROM airlines;
+-----+
| COUNT(*) |
+-----+
|      1000 |
+-----+
1 row in set (0.01 sec)

mysql> SELECT COUNT(*) FROM AirlineRating;
+-----+
| COUNT(*) |
+-----+
|      1000 |
+-----+
1 row in set (0.00 sec)
```

Connection Screenshot

Screenshot:

The screenshot shows a terminal window with the following content:

```
Use "gcloud config set project [PROJECT_ID]" to change to a different project.
m1_muktal13@cloudshell:~ (cs411-team058-helloworld)$ gcloud sql connect helloworld-058 --user=root
Allowlisting your IP for incoming connection for 5 minutes...done.
Connecting to database with SQL user [root].Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 22
Server version: 8.0.31-google (Google)

Copyright (c) 2000, 2023, Oracle and/or its affiliates.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> show databases;
+-----+
| Database      |
+-----+
| classimodels |
| flightshield |
| information_schema |
| mysql         |
| performance_schema |
| sys           |
+-----+
6 rows in set (0.01 sec)

mysql> use flightshield;
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
mysql> 
```

Advanced Queries:

Note: In this section, the value of delay is in minutes. A negative value of $-x$ indicates that the airline was **early** by x minutes. A positive value of y indicates that the airline was **delayed** by y minutes.

- 1) We are calculating the total average departure delay for every flight belonging to distinct airlines. We are joining the values from delays and airlines tables, and grouping by Airline name to calculate the aggregate average per airline.

```
SELECT a.AIRLINE, AVG(d.DEPARTURE_DELAY) AS Avg_Delay
FROM delays d
JOIN airlines a ON d.AIRLINE = a.IATA_CODE
GROUP BY a.AIRLINE
ORDER BY Avg_Delay DESC;
```

```
mysql> SELECT a.AIRLINE, AVG(d.DEPARTURE_DELAY) AS Avg_Delay
-> FROM delays d
-> JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND d.FLIGHT_NUMBER = f.FLIGHT_NUMBER JOIN airlines a ON a.IATA_CODE = f.AIRLINE
-> GROUP BY a.AIRLINE
-> ORDER BY Avg_Delay DESC;

+-----+-----+
| AIRLINE | Avg_Delay |
+-----+-----+
| Frontier Airlines Inc. | 19.5482 |
| American Eagle Airlines Inc. | 14.2178 |
| Spirit Air Lines | 13.6022 |
| Skywest Airlines Inc. | 12.1546 |
| United Air Lines Inc. | 12.1373 |
| JetBlue Airways | 10.6637 |
| American Airlines Inc. | 9.8684 |
| Atlantic Southeast Airlines | 9.2606 |
| Delta Air Lines Inc. | 6.8077 |
| Southwest Airlines Co. | 6.5269 |
| Virgin America | 6.4416 |
| US Airways Inc. | 5.2704 |
| Alaska Airlines Inc. | 3.1394 |
| Hawaiian Airlines Inc. | 1.0178 |
+-----+-----+
14 rows in set (2.02 sec)
```

The output is just 14 rows because our original data set only has information about 14 airlines (the remaining dummy airlines data we added does not have corresponding entries in the flights and delays table).

2) We are selecting the top 15 delayed flights due to DEPARTURE DELAYS based on the 4 parameters: Origin, Destination, Depart Date, and Return Date. This query simulates how a user might want to search for delayed flights based on their actual itinerary in our end application.

```
SELECT a.AIRLINE, AVG(d.DEPARTURE_DELAY) AS Avg_Delay
FROM delays d
JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND
d.FLIGHT_NUMBER = f.FLIGHT_NUMBER
JOIN airlines a ON a.IATA_CODE = f.AIRLINE
WHERE f.ORIGIN_AIRPORT = 'JFK'
AND f.DESTINATION_AIRPORT = 'SFO'
AND f.DATE = '1-Jan'
GROUP BY a.AIRLINE
ORDER BY Avg_Delay ASC;
```

```
mysql> SELECT a.AIRLINE, AVG(d.DEPARTURE_DELAY) AS Avg_Delay FROM delays d JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND d.FLIGHT_NUMBER = f.FLIGHT_NUMBER
JOIN airlines a ON a.IATA_CODE = f.AIRLINE WHERE f.ORIGIN_AIRPORT = 'JFK' AND f.DESTINATION_AIRPORT = 'SFO' AND f.DATE = '1-Jan' GROUP BY a.AIRLINE ORDER BY Avg_Delay ASC;
+-----+-----+
| AIRLINE | Avg_Delay |
+-----+-----+
| American Airlines Inc. | -5.7500 |
| JetBlue Airways | -1.6667 |
| United Air Lines Inc. | -0.6000 |
| Delta Air Lines Inc. | -0.2500 |
| Virgin America | 12.3333 |
+-----+-----+
5 rows in set (0.00 sec)
```

The output has just 5 rows because there are only 5 distinct airlines that serve the JFK to SFO route in our dataset.

3) Finally, similar to queries #2, we are calculating the average delay based on the total delay (i.e. departure delay plus arrival delay). We also simulated a real-world use case where a customer might want to check the anticipated delay for each airline, close to their holiday bookings. We tried superbowl weekend (30 Jan) and Valentine's day (14 Feb). We were able to assert that the delay values were significantly larger on these specific dates as anticipated.

```
SELECT a.AIRLINE, AVG(d.ARRIVAL_DELAY + d.DEPARTURE_DELAY) AS Avg_Delay
FROM delays d
JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND
d.FLIGHT_NUMBER = f.FLIGHT_NUMBER
JOIN airlines a ON a.IATA_CODE = f.AIRLINE
WHERE f.ORIGIN_AIRPORT = 'JFK'
AND f.DESTINATION_AIRPORT = 'PHX'
AND f.DATE = '30-JAN'
GROUP BY a.AIRLINE;
```

```
mysql> SELECT a.AIRLINE, AVG(d.ARRIVAL_DELAY + d.DEPARTURE_DELAY) AS Avg_Delay
-> FROM delays d
-> JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND d.FLIGHT_NUMBER = f.FLIGHT_NUMBER
-> JOIN airlines a ON a.IATA_CODE = f.AIRLINE
-> WHERE f.ORIGIN_AIRPORT = 'JFK'
-> AND f.DESTINATION_AIRPORT = 'PHX'
-> AND f.DATE = '30-JAN'
-> GROUP BY a.AIRLINE;
+-----+-----+
| AIRLINE | Avg_Delay |
+-----+-----+
| JetBlue Airways | 68.5000 |
| Delta Air Lines Inc. | 182.0000 |
| US Airways Inc. | 140.0000 |
+-----+-----+
3 rows in set (0.02 sec)
```

```
mysql> SELECT a.AIRLINE, AVG(d.ARRIVAL_DELAY + d.DEPARTURE_DELAY) AS Avg_Delay
-> FROM delays d
-> JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND d.FLIGHT_NUMBER = f.FLIGHT_NUMBER
-> JOIN airlines a ON a.IATA_CODE = f.AIRLINE
-> WHERE f.ORIGIN_AIRPORT = 'JFK'
-> AND f.DESTINATION_AIRPORT = 'SFO'
-> AND f.DATE = '14-Feb'
-> GROUP BY a.AIRLINE;
+-----+-----+
| AIRLINE | Avg_Delay |
+-----+-----+
| American Airlines Inc. | 127.2000 |
| JetBlue Airways | -2.5000 |
| Delta Air Lines Inc. | 60.2500 |
| United Air Lines Inc. | 37.8000 |
| Virgin America | -7.5000 |
+-----+-----+
5 rows in set (0.03 sec)
```

Similar to the output of query #2, this query also has just a few rows returned because there are only 5 distinct airlines that serve the JFK to SFO route in our dataset, and 3 distinct airlines that serve JFK to PHX route.

Part 2: Explain-Analyze and Indexing

Query #1:

```
SELECT a.AIRLINE, AVG(d.DEPARTURE_DELAY) AS Avg_Delay
FROM delays d
JOIN airlines a ON d.AIRLINE = a.IATA_CODE
GROUP BY a.AIRLINE
ORDER BY Avg_Delay DESC;
```

EXPLAIN ANALYZE without indexing:

```
| +-----+
| -> Sort: Avg_Delay DESC (actual time=29308.222..29308.223 rows=14 loops=1)
|   -> Table scan on <temporary> (actual time=29308.173..29308.181 rows=14 loops=1)
|     -> Aggregate using temporary table (actual time=29308.165..29308.165 rows=14 loops=1)
|       -> Nested loop inner join (cost=521909.75 rows=651275) (actual time=0.148..15539.678 rows=749794 loops=1)
|         -> Nested loop inner join (cost=293963.50 rows=651275) (actual time=0.136..4913.125 rows=749794 loops=1)
|           -> Covering index scan on f using AIRLINE (cost=66017.25 rows=651275) (actual time=0.073..1125.515 rows=749794 loops=1)
|             -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.25 rows=1) (actual time=0..004..0.004 rows=1 loops=749794)
|               -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.25 rows=1) (actual time=0.002..0.002 rows=1 loops=749794)
|             -> Filter: (d.AIRLINE = f.AIRLINE) (cost=0.25 rows=1) (actual time=0.014..0.014 rows=1 loops=749794)
|               -> Single-row index lookup on d using PRIMARY (DATE=f.DATE', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.25 rows=1) (actual time=0.013..0.013 rows=1 loops=749794)
|
+-----+
1 row in set (29.41 sec)
```

The original query is found to have a cost of **0.25** with a timing of **0.013s**.

- First we created an index based on the Departure Delays

Create index:

```
create index delay_idx on delays(DEPARTURE_DELAY);
```

```
| EXPLAIN
+
+-----+
|   > Sort: Avg_Delay DESC (actual_time=6900.015..6900.016 rows=14 loops=1)
|     -> Table scan on <temporary> (actual_time=6899.981..6899.986 rows=14 loops=1)
|       -> Aggregate using temporary table (actual_time=6899.975..6899.975 rows=14 loops=1)
|         -> Nested loop inner join (cost=521909.75 rows=651275) (actual_time=6.803..4694.694 rows=749794 loops=1)
|           -> Nested loop inner join (cost=293963.50 rows=651275) (actual_time=6.739..2128.531 rows=749794 loops=1)
|             -> Covering index scan on f using AIRLINE (cost=66017.25 rows=651275) (actual_time=6.687..1479.672 rows=749794 loops=1)
|               -> Single-row Index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.25 rows=1) (actual_time=0.001..0.001 rows=1 loops=749794)
|                 -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.25 rows=1) (actual_time=0.001..0.001 rows=1 loops=749794)
|                   -> Filter: (d.AIRLINE = f.AIRLINE) (cost=0.25 rows=1) (actual_time=0.003..0.003 rows=1 loops=749794)
|                     -> Single-row index lookup on d using PRIMARY (DATE=f.DATE, AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.25 rows=1) (actual_time=0.003..0.003
rows=1 loops=749794)
|
+-----+
1 row in set (6.91 sec)

mysql>
```

The indexed query on delay_idx on flights is found to have an unchanged cost of **0.25** with a timing of **0.003s**. Hence while there is an improvement in timing, the cost was unchanged, hence we continued to try other options.

- Next, we attempted to index based on Airline names.

```
create index airlines_idx on airlines(AIRLINE(5));
```

```
| EXPLAIN
+
+-----+
|   > Sort: Avg_Delay DESC (actual_time=5836.630..5836.631 rows=14 loops=1)
|     -> Table scan on <temporary> (actual_time=5836.596..5836.600 rows=14 loops=1)
|       -> Aggregate using temporary table (actual_time=5836.589..5836.589 rows=14 loops=1)
|         -> Nested loop inner join (cost=521909.75 rows=651275) (actual_time=0.090..3357.508 rows=749794 loops=1)
|           -> Nested loop inner join (cost=293963.50 rows=651275) (actual_time=0.078..974.832 rows=749794 loops=1)
|             -> Covering index scan on f using AIRLINE (cost=66017.25 rows=651275) (actual_time=0.060..285.987 rows=749794 loops=1)
|               -> Single-row Index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.25 rows=1) (actual_time=0.001..0.001 rows=1 loops=749794)
|                 -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.25 rows=1) (actual_time=0.001..0.001 rows=1 loops=749794)
|                   -> Filter: (d.AIRLINE = f.AIRLINE) (cost=0.25 rows=1) (actual_time=0.003..0.003 rows=1 loops=749794)
|                     -> Single-row index lookup on d using PRIMARY (DATE=f.DATE, AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.25 rows=1) (actual_time=0.003..0.003
rows=1 loops=749794)
|
+-----+
1 row in set (5.84 sec)
```

For the indexed query on Airline Names, there is found to be an unchanged cost of **0.25** with a timing of **0.003s**. The reason for the cost remaining unchanged and the time improving is probably because the database only has 14 airlines (the remaining dummy airlines data we added does not have corresponding entries in the flights and delays table).**The second reason is that since airline is a primary key, it is already indexed.** Hence as there are no cost improvements, it is determined that indexing based on Airlines is not an optimal indexing.

- Then, we attempted to index based on both Airline names and Departure Delays.

```
create index airlines_idx on airlines(AIRLINE(5));
create index delay_idx on delays(DEPARTURE_DELAY);
```

```
--+
| EXPLAIN
+-----+
|   |   |
|   +-- Sort: Avg Delay DESC (actual time=5125.312..5125.313 rows=14 loops=1)
|   |   +-- Table scan on <temporary> (actual time=5125.279..5125.284 rows=14 loops=1)
|   |       +-- Aggregate using temporary table (actual time=5125.273..5125.273 rows=14 loops=1)
|   |           +-- Nested loop inner join (cost=521909.75 rows=651275) (actual time=0.086..3026.006 rows=749794 loops=1)
|   |               +-- Covering index scan on f using AIRLINE (cost=66017.25 rows=651275) (actual time=0.052..253.896 rows=749794 loops=1)
|   |                   +-- Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.25 rows=1) (actual time=0.001..0.001 rows=1 loops=749794)
|   |                       +-- Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=749794)
|   |                           +-- Filter: (d.AIRLINE = f.AIRLINE) (cost=0.25 rows=1) (actual time=0.003..0.003 rows=1 loops=749794)
|   |                               +-- Single-row index lookup on d using PRIMARY (DATE=f.DATE , AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.25 rows=1) (actual time=0.002..0.002
|   |                                   rows=1 loops=749794)
|   |
|   +-----+
1 row in set (5.12 sec)
```

For the indexed query on both delay_idx and airlines_idx, there is found to be an unchanged cost of 0.25 with a timing improvement of 0.002s. Which can be attributed to the database only having 14 airlines as seen in the earlier indexing and airline being the primary key is already indexed.

- Then, we attempted to index based on Date, Airline names and Flight Number.

```
create index delay_idx on delays(DATE(5), AIRLINE(3),
FLIGHT_NUMBER);
```

```
1 row in set (1.30 sec)

mysql> create index delay_idx on delays(DATE(5), AIRLINE(3),FLIGHT_NUMBER);
Query OK, 0 rows affected (6.44 sec)
Records: 0  Duplicates: 0  Warnings: 0

mysql> EXPLAIN ANALYZE SELECT a.AIRLINE, AVG(d.DEPARTURE_DELAY) AS Avg_Delay FROM delays d JOIN airlines a ON d.AIRLINE = a.IATA_CODE GROUP BY a.AIRLINE ORDER BY Avg_Delay DESC;
+-----+
| EXPLAIN
+-----+
|   |
+-----+
| -> Sort: Avg Delay DESC (actual time=1512.044..1512.045 rows=14 loops=1)
    -> Table scan on temporary table (actual time=1512.011..1512.016 rows=14 loops=1)
        -> Aggregate using temporary table (actual time=1512.008..1512.008 rows=14 loops=1)
            -> Nested loop inner join (cost=350460.60 rows=774078) (actual time=0.189..892.643 rows=749794 loops=1)
                -> Table scan on d  (cost=79533.30 rows=774078) (actual time=0.169..328.142 rows=749794 loops=1)
                    -> Filter: (d.AIRLINE = a.IATA_CODE) (cost=0.25 rows=1) (actual time=0.000..0.001 rows=1 loops=749794)
                        -> Single-row index lookup on a using PRIMARY (IATA_CODE=d.AIRLINE) (cost=0.25 rows=1) (actual time=0.000..0.000 rows=1 loops=749794)
|
+-----+
|   |
+-----+
1 row in set (1.52 sec)
```

For the indexed query on both delay_idx and airlines_idx, there is found to be an unchanged cost of **0.25** with a timing improvement to **0.000s**. Since all the three attributes are primary keys, Indexing has no effect.

Hence we decided to go with the non-indexed or default query, since adding an index does not seem to be very helpful in this situation, and is not reducing cost and timing both in any instances.

Query #2:

```

SELECT a.AIRLINE, AVG(d.DEPARTURE_DELAY) AS Avg_Delay
FROM delays d
JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND
d.FLIGHT_NUMBER = f.FLIGHT_NUMBER
JOIN airlines a ON a.IATA_CODE = f.AIRLINE
WHERE f.ORIGIN_AIRPORT = 'JFK'
AND f.DESTINATION_AIRPORT = 'SFO'
AND f.DATE = '1-Jan'
GROUP BY a.AIRLINE
ORDER BY Avg_Delay ASC;

```

Explain Analyze without Indexing:

```

+-----+
| -> Sort: Avg_Delay (actual time=1.279..1.279 rows=5 loops=1)
   -> Table scan on <temporary> (actual time=1.249..1.250 rows=5 loops=1)
      -> Aggregate using temporary table (actual time=1.246..1.246 rows=5 loops=1)
         -> Nested loop inner join (cost=57.97 rows=1) (actual time=0.193..1.188 rows=19 loops=1)
            -> Nested loop inner join (cost=57.62 rows=1) (actual time=0.168..1.010 rows=19 loops=1)
               -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') and (f.DESTINATION_AIRPORT = 'SFO')) and (f.DATE = '1-Jan')) (cost=57.27 rows=1) (actual time=0.152..0.959 rows=19 loops=1)
s=1)
               -> Intersect rows sorted by row ID (cost=57.27 rows=1) (actual time=0.088..0.883 rows=19 loops=1)
                  -> Index range scan on f using ORIGIN_AIRPORT over (ORIGIN_AIRPORT = 'JFK' AND DATE = '1-Jan') (cost=25.51 rows=272) (actual time=0.034..0.382 rows=272 loops=1)
                     -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'SFO' AND DATE = '1-Jan') (cost=31.66 rows=341) (actual time=0.049..0.51 rows=294 loops=1)
                     -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
                        -> Single-row Index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
                           -> Filter: ((d.AIRLINE = f.AIRLINE) and (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.009..0.009 rows=1 loops=19)
                           -> Single-row index lookup on d using PRIMARY (DATE='1-Jan', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.008..0.008 rows=1 loops=19)
rows=1 loops=19)
|
+-----+
1 row in set (0.01 sec)

```

The original query is found to have a cost of 0.35 with a timing of 0.008s.

- Since the query calculates the average departure delay, an index on this column could potentially improve the speed of the aggregation

```

CREATE INDEX idx_delays_departure_delay ON delays
(DEPARTURE_DELAY);

```

```
| EXPLAIN
+
+-----+
|   +-- Sort: Avg_Delay (actual time=0.756..0.756 rows=5 loops=1)
|   +-- Table scan on <temporary> (actual time=0.734..0.735 rows=5 loops=1)
|       +-- Aggregate using temporary table (actual time=0.733..0.733 rows=5 loops=1)
|           +-- Nested loop inner join (cost=58.18 rows=1) (actual time=0.094..0.661 rows=19 loops=1)
|               +-- Nested loop inner join (cost=57.83 rows=1) (actual time=0.083..0.548 rows=19 loops=1)
|                   +-- Filter: (f.ORIGIN_AIRPORT = 'JFK' AND f.DESTINATION_AIRPORT = 'SFO' AND f.DATE = '1-Jan') (cost=57.48 rows=1) (actual time=0.069..0.504 rows=19 loops=1)
|                       +-- Index range scan on f using origin_idx over (ORIGIN_AIRPORT = 'JFK' AND DATE = '1-Jan') (cost=57.48 rows=19 loops=1)
|                           +-- Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|                               +-- Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.001..0.002 rows=1 loops=19)
|                                   +-- Filter: (d.DATE = f.DATE) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=19)
|                                       +-- Single-row index lookup on d using PRIMARY (DATE='1-Jan', AIRLINE=f.AIRLINE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=19)
|                                           +-- Single-row index lookup on d using PRIMARY (DATE='1-Jan', AIRLINE=f.AIRLINE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=19)
```

The indexed query on flights (origin_airport , destination_airport) is found to have an unchanged cost of **0.35** with a timing of **0.006s**. Hence as the cost was unchanged with an insignificant improvement in timing, we continued to try other options.

- Since the query filters results based on specific origin and destination airports, an index on these columns would help the database quickly find relevant rows.

```
CREATE      INDEX      idx_flights_origin_destination      ON      flights
(ORIGIN_AIRPORT, DESTINATION_AIRPORT);
```

```
| EXPLAIN
+
+-----+
|   +-- Sort: Avg_Delay (actual time=0.732..0.733 rows=5 loops=1)
|   +-- Table scan on <temporary> (actual time=0.710..0.712 rows=5 loops=1)
|       +-- Aggregate using temporary table (actual time=0.709..0.709 rows=5 loops=1)
|           +-- Nested loop inner join (cost=58.18 rows=1) (actual time=0.095..0.664 rows=19 loops=1)
|               +-- Nested loop inner join (cost=57.83 rows=1) (actual time=0.083..0.548 rows=19 loops=1)
|                   +-- Filter: (f.ORIGIN_AIRPORT = 'JFK' AND f.DESTINATION_AIRPORT = 'SFO' AND f.DATE = '1-Jan') (cost=57.48 rows=1) (actual time=0.070..0.511 rows=19 loops=1)
|                       +-- Index range scan on f using origin_idx over (ORIGIN_AIRPORT = 'JFK' AND DATE = '1-Jan') (cost=57.48 rows=19 loops=1)
|                           +-- Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.067..0.502 rows=19 loops=1)
|                               +-- Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'SFO' AND DATE = '1-Jan') (cost=25.57 rows=272) (actual time=0.028..0.233 rows=272 loops=1)
|                                   +-- Filter: (d.DATE = f.DATE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|                                       +-- Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.001..0.001 rows=1 loops=19)
|                                           +-- Filter: (d.DATE = f.DATE) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=19)
|                                               +-- Single-row index lookup on d using PRIMARY (DATE='1-Jan', AIRLINE=f.AIRLINE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=19)
```

The indexed query on flights (origin_airport , destination_airport) is found to have an unchanged cost of **0.35** with a timing of **0.006s**. Hence as the cost was unchanged with an insignificant improvement in timing, we continued to try other options.

- As there are many dates first we tried to create an index on the date attribute and compare it to the original query:

```
create index date_idx on flights(DATE(5));
```

```

-----+
| -> Sort: Avg_Delay (actual time=0.847..0.848 rows=5 loops=1)
|   -> Table scan on <temporary> (actual time=0.827..0.828 rows=5 loops=1)
|     -> Aggregate using temporary table (actual time=0.825..0.825 rows=5 loops=1)
|       -> Nested loop inner join (cost=58.35 rows=1) (actual time=0.142..0.777 rows=19 loops=1)
|         -> Nested loop inner join (cost=58.00 rows=1) (actual time=0.128..0.630 rows=19 loops=1)
|           -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') and (f.DESTINATION_AIRPORT = 'SFO')) and (f.DATE = '1-Jan')) (cost=57.65 rows=1) (actual time=0.105..0.578 rows=19 loops=1)
|             s=1
|               -> Intersect rows sorted by row ID (cost=57.65 rows=1) (actual time=0.102..0.568 rows=19 loops=1)
|                 -> Index range scan on f using ORIGIN_AIRPORT over (ORIGIN_AIRPORT = 'JFK' AND DATE = '1-Jan') (cost=25.64 rows=272) (actual time=0.036..0.258 rows=272
|                   loops=1)
|                   -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'SFO' AND DATE = '1-Jan') (cost=31.91 rows=341) (actual time=0.062..0.24
|                   2 rows=294 loops=1)
|                     -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|                       -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|                         -> Filter: (d.DATE = f.DATE) (cost=0.35 rows=1) (actual time=0.007..0.007 rows=1 loops=19)
|                           -> Single-row index lookup on d using PRIMARY (DATE='1-Jan', AIRLINE=f.AIRLINE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.007..0.007 row
|                           s=1 loops=19)
|             |
+-----+

```

The indexed query on flight Dates is found to have an unchanged cost of **0.35** with a timing of **0.007s**. Hence as the cost was unchanged with an insignificant improvement in timing, this is because Date is already indexed due to it being the primary key.

- We tried to create an index on the Origin Airport attribute.

```
create index origin_idx on flights(ORIGIN_AIRPORT(3));
```

```

-----+
| -> Sort: Avg_Delay (actual time=0.677..0.677 rows=5 loops=1)
|   -> Table scan on <temporary> (actual time=0.659..0.660 rows=5 loops=1)
|     -> Aggregate using temporary table (actual time=0.658..0.658 rows=5 loops=1)
|       -> Nested loop inner join (cost=58.28 rows=1) (actual time=0.091..0.615 rows=19 loops=1)
|         -> Nested loop inner join (cost=57.93 rows=1) (actual time=0.082..0.514 rows=19 loops=1)
|           -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') and (f.DESTINATION_AIRPORT = 'SFO')) and (f.DATE = '1-Jan')) (cost=57.58 rows=1) (actual time=0.068..0.473 rows=19 loops=1)
|             s=1
|               -> Intersect rows sorted by row ID (cost=57.58 rows=1) (actual time=0.065..0.464 rows=19 loops=1)
|                 -> Index range scan on f using ORIGIN_AIRPORT over (ORIGIN_AIRPORT = 'JFK' AND DATE = '1-Jan') (cost=25.60 rows=272) (actual time=0.030..0.195 rows=272
|                   loops=1)
|                   -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'SFO' AND DATE = '1-Jan') (cost=31.88 rows=341) (actual time=0.030..0.20
|                   5 rows=294 loops=1)
|                     -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|                       -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|                         -> Filter: (d.DATE = f.DATE) (cost=0.35 rows=1) (actual time=0.005..0.005 rows=1 loops=19)
|                           -> Single-row index lookup on d using PRIMARY (DATE='1-Jan', AIRLINE=f.AIRLINE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.005..0.005 row
|                           s=1 loops=19)
|             |
+-----+

```

The indexed query on Origin Airports was also found to have an unchanged cost of **0.35** with a timing of **0.005s**. Thus the cost is unchanged with an insignificant improvement in timing, we try another option.

- This time we combined both the above attributes, together with destination, to see if that would optimize performance.

```
create index all_idx on flights(ORIGIN_AIRPORT(3),
                               DESTINATION_AIRPORT(3), DATE(5));
```

```
| -> Sort: Avg_Delay (actual time=0.284..0.285 rows=5 loops=1)
|   -> Table scan on <temporary> (actual time=0.230..0.231 rows=5 loops=1)
|     -> Aggregate using temporary table (actual time=0.228..0.228 rows=5 loops=1)
|       -> Nested loop inner join (cost=19.95 rows=19) (actual time=0.050..0.166 rows=19 loops=1)
|         -> Nested loop inner join (cost=13.30 rows=19) (actual time=0.042..0.073 rows=19 loops=1)
|           -> Index lookup on f using all_idx (ORIGIN_AIRPORT='JFK', DESTINATION_AIRPORT='SFO', DATE='1-Jan', DATE='1-Jan') (cost=6.65 rows=19) (actual time=0.025..0.031
rows=19 loops=1)
|             -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.26 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|               -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.26 rows=1) (actual time=0.002..0.002 rows=1 loops=19)
|                 -> Filter: (d.DATE = f.DATE) (cost=0.26 rows=1) (actual time=0.005..0.005 rows=1 loops=19)
|                   -> Single-row index lookup on d using PRIMARY (DATE='1-Jan', AIRLINE=f.AIRLINE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.26 rows=1) (actual time=0.004..0.004 row
s=1 loops=19)
|
+-----+
1 row in set (0.00 sec)
```

The indexed query on the Origin and Destination airports and the Date was found to also have an improved cost of **0.26** with a timing of **0.004s**. Thus the cost reduced significantly with an improvement in timing, as it's faster to pinpoint exact flights. We choose this indexing for this query due to significant improvement in the cost and time.

Query #3:

```

SELECT a.AIRLINE, AVG(d.ARRIVAL_DELAY + d.DEPARTURE_DELAY) AS
Avg_Delay
FROM delays d
JOIN flights f ON d.DATE = f.DATE AND d.AIRLINE = f.AIRLINE AND
d.FLIGHT_NUMBER = f.FLIGHT_NUMBER
JOIN airlines a ON a.IATA_CODE = f.AIRLINE
WHERE f.ORIGIN_AIRPORT = 'JFK'
AND f.DESTINATION_AIRPORT = 'PHX'
AND f.DATE = '30-JAN'
GROUP BY a.AIRLINE;

```

EXPLAIN ANALYZE without indexing:

```

-----+
| > Table scan on <temporary> (actual time=0.593..0.594 rows=3 loops=1)
|   -> Aggregate using temporary Table (actual time=0.592..0.592 rows=3 loops=1)
|     -> Nested loop inner join (cost=60.13 rows=1) (actual time=0.141..0.561 rows=8 loops=1)
|       -> Nested loop inner join (cost=59.78 rows=1) (actual time=0.129..0.508 rows=8 loops=1)
|         -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') AND (f.DESTINATION_AIRPORT = 'PHX')) (cost=58.93 rows=1) (actual time=0.116..0.483 rows=8 loops=1)
|           -> Intersect rows sorted by row ID (cost=58.93 rows=1) (actual time=0.113..0.477 rows=8 loops=1)
|             -> Index range scan on f using origin_idx over (ORIGIN_AIRPORT = 'JFK' AND DATE = '30-JAN') (cost=26.24 rows=283) (actual time=0.027..0.226 rows=283 loops=1)
|               -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'PHX' AND DATE = '30-JAN') (cost=32.59 rows=349) (actual time=0.028..0.187 rows=269 loops=1)
|                 -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.85 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
|                   -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.85 rows=1) (actual time=0.002..0.002 rows=1 loops=8)
|                     -> Filter: ((d.AIRLINE = f.AIRLINE) AND (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
|                       -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
|                         -> Filter: (d.ARRIVAL_DELAY + d.DEPARTURE_DELAY) > 0 (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
|                           -> Aggregate using temporary Table (actual time=0.006..0.006 rows=1 loops=1)
|                             -> Final Aggregate (actual time=0.006..0.006 rows=1 loops=1)
|                               -> Result (actual time=0.006..0.006 rows=1 loops=1)
-----+

```

The original query is found to have a cost of 0.35 with a timing of 0.006s.

With indexing:

- For this query we tried to create an index on the Arrival Delays and compare it to the original query:

```
create index average_idx on delays(ARRIVAL_DELAY);
```

```

| --> Table scan on <temporary> (actual time=0.560..0.561 rows=3 loops=1)
  -> Aggregate using temporary table (actual time=0.559..0.559 rows=3 loops=1)
    -> Nested loop inner join (cost=59.63 rows=1) (actual time=0.143..0.527 rows=8 loops=1)
      -> Nested loop inner join (cost=59.28 rows=1) (actual time=0.130..0.474 rows=8 loops=1)
        -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') and (f.DESTINATION_AIRPORT = 'PHX')) (cost=58.93 rows=1) (actual time=0.118..0.450 rows=8 loops=1)
          -> Intersect rows sorted by row ID (cost=58.93 rows=1) (actual time=0.114..0.444 rows=8 loops=1)
            -> Index range scan on f using origin_idx over (ORIGIN_AIRPORT = 'JFK' AND DATE = '30-JAN') (cost=26.24 rows=283) (actual time=0.024..0.186 rows=283 loops=1)
1)   -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'PHX' AND DATE = '30-JAN') (cost=32.59 rows=349) (actual time=0.032..0.192 r
ows=269 loops=1)
      -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
        -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=8)
      -> Filter: ((d.AIRLINE = f.AIRLINE) and (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
        -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows
=1 loops=8)
+
+-----+

```

The indexed query on Arrival Delays was found to have an unchanged cost of **0.35** with a timing of **0.006s**. Thus as the cost is unchanged with an insignificant change in timing, we try another option.

- We tried to create an index on the Departure Delays attribute.

```
create index dept_idx on delays(DEPARTURE_DELAY);
```

```

| --> Table scan on <temporary> (actual time=0.582..0.583 rows=3 loops=1)
  -> Aggregate using temporary table (actual time=0.581..0.581 rows=3 loops=1)
    -> Nested loop inner join (cost=59.63 rows=1) (actual time=0.146..0.548 rows=8 loops=1)
      -> Nested loop inner join (cost=59.28 rows=1) (actual time=0.132..0.494 rows=8 loops=1)
        -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') and (f.DESTINATION_AIRPORT = 'PHX') and (f.DATE = '30-JAN')) (cost=58.93 rows=1) (actual time=0.120..0.470 rows=8 loops=1)
          -> Intersect rows sorted by row ID (cost=58.93 rows=1) (actual time=0.117..0.464 rows=8 loops=1)
            -> Index range scan on f using origin_idx over (ORIGIN_AIRPORT = 'JFK' AND DATE = '30-JAN') (cost=26.24 rows=283) (actual time=0.029..0.194 rows=283 loops=1)
1)   -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'PHX' AND DATE = '30-JAN') (cost=32.59 rows=349) (actual time=0.030..0.203 r
ows=269 loops=1)
      -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
        -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=8)
      -> Filter: ((d.AIRLINE = f.AIRLINE) and (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
        -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows
=1 loops=8)
+
+-----+

```

The indexed query on Departure Delays was found to have an unchanged cost of **0.35** with a timing of **0.006s**. Thus as the cost is unchanged with no change in timing, we try another option.

- We tried to create an index using both the Departure and Arrival Delays attribute.

```
create index avg_idx on delays(DEPARTURE_DELAY, ARRIVAL_DELAY);
```

```

| -> Table scan on <temporary> (actual time=0.620..0.621 rows=3 loops=1)
  -> Aggregate using temporary table (actual time=0.619..0.619 rows=3 loops=1)
    -> Nested loop inner join (cost=59.63 rows=1) (actual time=0.144..0.579 rows=8 loops=1)
      -> Nested loop inner join (cost=59.28 rows=1) (actual time=0.131..0.526 rows=8 loops=1)
        -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') and (f.DESTINATION_AIRPORT = 'PHX') and (f.DATE = '30-JAN')) (cost=58.93 rows=1) (actual time=0.120..0.502 rows=8 loops=1)
          -> Intersect rows sorted by row ID (cost=58.93 rows=1) (actual time=0.117..0.496 rows=8 loops=1)
            -> Index range scan on f using origin_idx over (ORIGIN_AIRPORT = 'JFK' AND DATE = '30-JAN') (cost=26.24 rows=283) (actual time=0.026..0.205 rows=283 loops=1)
              -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'PHX' AND DATE = '30-JAN') (cost=32.59 rows=349) (actual time=0.033..0.195 rows=269 loops=1)
                -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.003 rows=1 loops=8)
                  -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.002..0.002 rows=1 loops=8)
                    -> Filter: ((d.AIRLINE = f.AIRLINE) and (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
                      -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
+-----+
1 row in set (0.01 sec)

```

The indexed query using both Departure and Arrival Delays was also found to have an unchanged cost of **0.35** with a timing of **0.006s**. Thus as the cost is unchanged with no change in timing, we try another option.

- Hence we tried to create an index combining both the Departure and Arrival Delays attribute.

```
create index avg_idx on delays(DEPARTURE_DELAY+ARRIVAL_DELAY);
```

```

| -> Table scan on <temporary> (actual time=0.651..0.652 rows=3 loops=1)
| -> Aggregate using temporary table (actual time=0.643..0.643 rows=3 loops=1)
|   -> Nested loop inner join (cost=59.63 rows=1) (actual time=0.175..0.558 rows=8 loops=1)
|     -> Nested loop inner join (cost=59.28 rows=1) (actual time=0.161..0.504 rows=8 loops=1)
|       -> Filter: ((f.ORIGIN_AIRPORT = 'JFK') AND (f.DESTINATION_AIRPORT = 'PHX') AND (f.DATE = '30-JAN')) (cost=58.93 rows=1) (actual time=0.129..0.459 rows=8 loops=1)
|         -> Intersect rows sorted by row ID (cost=58.93 rows=1) (actual time=0.123..0.450 rows=8 loops=1)
|           -> Index range scan on f using origin_idx over (ORIGIN_AIRPORT = 'JFK' AND DATE = '30-JAN') (cost=26.24 rows=283) (actual time=0.029..0.190 rows=283 loops=1)
|             -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'PHX' AND DATE = '30-JAN') (cost=32.59 rows=349) (actual time=0.036..0.194 rows=269 loops=1)
|               -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.005..0.005 rows=1 loops=8)
|                 -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.004..0.004 rows=1 loops=8)
|                   -> Filter: ((d.AIRLINE = f.AIRLINE) AND (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
|                     -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
|
+

```

The indexed query combining both Departure and Arrival Delays was also found to have an unchanged cost of **0.35** with a timing of **0.006s**. Thus as the cost is unchanged with no change in timing, we try another option.

- We tried to create an index using both the Departure and Arrival Delays attribute as well as the Flight Dates.

```
CREATE INDEX avg_idx ON delays(DATE(5), DEPARTURE_DELAY,
ARRIVAL_DELAY);
```

```
|----+
| -> Table scan on <temporary> (actual time=0.768..0.770 rows=3 loops=1)
|   -> Aggregate using temporary table (actual time=0.767..0.767 rows=3 loops=1)
|     -> Nested loop inner join (cost=59.39 rows=1) (actual time=0.227..0.761 rows=3 loops=1)
|       -> Nested loop inner join (cost=59.28 rows=1) (actual time=0.209..0.670 rows=8 loops=1)
|         -> Filter: ((f.ORIGIN AIRPORT = 'JFK') AND (f.DESTINATION AIRPORT = 'PHX')) AND (f.DATE = '30-JAN')) (cost=58.93 rows=1) (actual time=0.194..0.641 rows=8 loops=1)
|           -> Intersect rows sorted by row ID (cost=58.93 rows=1) (actual time=0.190..0.634 rows=8 loops=1)
|             -> Index range scan on f using origin_idx over (ORIGIN AIRPORT = 'JFK' AND DATE = '30-JAN') (cost=26.24 rows=283) (actual time=0.025..0.291 rows=283 loops=1)
|               -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'PHX' AND DATE = '30-JAN') (cost=32.59 rows=349) (actual time=0.045..0.278 r
ows=269 loops=1)
|                 -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
|                   -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
|                     -> Filter: ((d.AIRLINE = f.AIRLINE) AND (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.007..0.007 rows=1 loops=8)
|                       -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.007..0.007 rows=1 loops=8)
|               |
|----+
```

The indexed query combining both Departure and Arrival Delays was also found to have an unchanged cost of **0.35** with an insignificantly changed time of **0.007s**. **Date is a primary key which is already Indexed**. Thus as the cost is unchanged with almost no change in timing, we try another option.

- We tried to create an index using the average of the Departure and Arrival Delays

```
create index avg_idx on delays(
((DEPARTURE_DELAY+ ARRIVAL_DELAY)/2));
```

```
|----+
| -> Table scan on <temporary> (actual time=0.672..0.673 rows=3 loops=1)
|   -> Aggregate using temporary table (actual time=0.670..0.670 rows=3 loops=1)
|     -> Nested loop inner join (cost=59.63 rows=1) (actual time=0.252..0.628 rows=8 loops=1)
|       -> Nested loop inner join (cost=59.52 rows=1) (actual time=0.225..0.552 rows=8 loops=1)
|         -> Filter: ((f.ORIGIN AIRPORT = 'JFK') AND (f.DESTINATION AIRPORT = 'PHX')) AND (f.DATE = '30-JAN')) (cost=58.93 rows=1) (actual time=0.209..0.533 rows=8 loops=1)
|             -> Intersect rows sorted by row ID (cost=58.93 rows=1) (actual time=0.204..0.526 rows=8 loops=1)
|               -> Index range scan on f using origin_idx over (ORIGIN AIRPORT = 'JFK' AND DATE = '30-JAN') (cost=26.24 rows=283) (actual time=0.040..0.203 rows=283 loops=1)
|                 -> Index range scan on f using DESTINATION_AIRPORT over (DESTINATION_AIRPORT = 'PHX' AND DATE = '30-JAN') (cost=32.59 rows=349) (actual time=0.104..0.256 r
ows=269 loops=1)
|                   -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.35 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
|                     -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.35 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
|                       -> Filter: ((d.AIRLINE = f.AIRLINE) AND (d.DATE = f.DATE)) (cost=0.35 rows=1) (actual time=0.008..0.008 rows=1 loops=8)
|                         -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.35 rows=1) (actual time=0.007..0.007 rows=1 loops=8)
|               |
|----+
```

The indexed query averaging the Departure and Arrival Delays was also found to have an unchanged cost of **0.35** with an insignificantly changed time of **0.007s**. Thus as the cost is unchanged with almost no change in timing, we try another option.

- We tried to **combine the indexes** using the Date, Departure and Arrival Delays, as well as the other index using Origin and Destination airports.

```
CREATE INDEX avg_idx ON delays(DATE(5), DEPARTURE_DELAY,
ARRIVAL_DELAY);
CREATE INDEX flg_idx ON flights(ORIGIN_AIRPORT(3),
DESTINATION_AIRPORT(3));
```

```
| EXPLAIN
+-----+
| > Table scan on <temporary> (actual time=0.183..0.184 rows=3 loops=1)
|   -> Aggregate using temporary table (actual time=0.181..0.181 rows=3 loops=1)
|     -> Nested loop inner join (cost=8.04 rows=8) (actual time=0.087..0.143 rows=8 loops=1)
|       -> Nested loop inner join (cost=5.24 rows=8) (actual time=0.071..0.088 rows=8 loops=1)
|         -> Covering index lookup on f using flg_idx (ORIGIN_AIRPORT='JFK', DESTINATION_AIRPORT='PHX', DATE='30-JAN') (cost=2.44 rows=8) (actual time=0.029..0.034 rows=8 loops=1)
|           -> Filter: (a.IATA_CODE = f.AIRLINE) (cost=0.26 rows=1) (actual time=0.003..0.003 rows=1 loops=8)
|             -> Single-row index lookup on a using PRIMARY (IATA_CODE=f.AIRLINE) (cost=0.26 rows=1) (actual time=0.002..0.003 rows=1 loops=8)
|               -> Filter: ((d.AIRLINE = f.AIRLINE) AND (d.DATE = f.DATE)) (cost=0.26 rows=1) (actual time=0.007..0.007 rows=1 loops=8)
|                 -> Single-row index lookup on d using PRIMARY (DATE='30-JAN', AIRLINE=a.IATA_CODE, FLIGHT_NUMBER=f.FLIGHT_NUMBER) (cost=0.26 rows=1) (actual time=0.006..0.006 rows=1 loops=8)
+-----+
1 row in set (0.01 sec)
```

The combined query finally resulted in a reduced cost of **0.26** with an insignificantly changed time of **0.006s**. Thus as the cost is reduced with slight change in timing, we treat this as the optimized query and the choice for selecting this as our index on this query.