

## Lab 3 Anna Ślęzak

### 1. Wykonaj zapytania (3p):

1. Sprawdź dla każdego dnia tygodnia w jakiej godzinie jest najwięcej branych taksówek.

```
1 with ranked_hours as (  
2     select day_of_week(tpep_pickup_datetime) as day_of_week,  
3     split_part(  
4         split_part(cast(tpep_pickup_datetime as varchar), ' ', 2),  
5         ':',  
6         1  
7     ) as hour,  
8     count(*) as no_of_rides,  
9     row_number() over (  
10        partition by day_of_week(tpep_pickup_datetime)  
11        order by count(*) desc  
12    ) as rank  
13 from TaxiDataYellow  
14 group by day_of_week(tpep_pickup_datetime),  
15 split_part(  
16     split_part(cast(tpep_pickup_datetime as varchar), ' ', 2),  
17     ':',  
18     1  
19 )  
20 )  
21 select day_of_week,  
22        hour,  
23        no_of_rides  
24 from ranked_hours  
25 where rank = 1;
```

2. Znajdź przejazdy które znacząco odbiegają od ceny standardowej (outliery)

```
with fare_stats as (  
    select avg(fare_amount) as avg_fare,  
           stddev(fare_amount) as stddev_fare  
    from taxidatayellow  
)  
select *  
from taxidatayellow  
    cross join fare_stats  
where fare_amount < avg_fare - 1.5 * stddev_fare  
    or fare_amount > avg_fare + 1.5 * stddev_fare;
```

3. Znajdź średnią cenę przejazdu per osoba dla każdej firmy w zależności od liczby osób (z przedziału 1-7) w samochodzie, zlicz też liczbę przejazdów dla każdej ceny razem z procentowym udziałem we wszystkich przejazdach danej firmy.

```

SELECT vendorid,
       passenger_count,
       AVG(fare_amount / passenger_count) AS avg_price_per_person,
       COUNT(*) AS ride_count,
       (COUNT(*) * 100.0) / SUM(COUNT(*)) OVER (PARTITION BY vendorid) AS percentage_of_total_rides
FROM TaxiDataYellow
WHERE passenger_count BETWEEN 1 AND 7
GROUP BY vendorid,
       passenger_count;

```

2. Znajdź własny dataset w formacie innym niż Parquet, skonwertuj go do formatu Parquet, a następnie wykonaj własne nietrywialne zapytanie na obu formatach z wykorzystaniem AWS Athena i sprawdź które zapytanie działa szybciej, które procesuje więcej danych. Sprawdź równoważność zwróconych wyników. Do konwersji wystarczą 2 zapytania.

(2p) <https://docs.aws.amazon.com/athena/latest/ug/ctas-examples.html#ctas-example-format>

Tabela w formacie csv:

```

CREATE EXTERNAL TABLE IF NOT EXISTS dataset_flight (
  airline STRING,
  date_of_journey STRING,
  source STRING,
  destination STRING,
  route STRING,
  dep_time STRING,
  arrival_time STRING,
  duration STRING,
  total_stops STRING,
  additional_info STRING,
  price bigint
)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION 's3://844124137610-us-east-1-athena-results-bucket-qhv015up15/dataset_flights/'

```

Tabela w formacie parquet:

```

1 CREATE TABLE flights_parquet
2 WITH (
3     format = 'Parquet',
4     write_compression = 'SNAPPY')
5 AS SELECT *
6 FROM dataset_flight;
7
8 select * from flights_parquet
9

```

**Wybór najpopularniejszej trasy w danym miesiącu dla określonej linii lotniczej oraz średnia cena biletu dla tej trasy.**

CSV:

```

1 WITH MonthlyAverages AS (
2     SELECT
3         airline,
4         split_part(date_of_journey, '/', 2) AS travel_month,
5         destination,
6         AVG(price) AS avg_price,
7         RANK() OVER (PARTITION BY airline, split_part(date_of_journey, '/', 2) ORDER BY COUNT(*) DESC) AS rank
8     FROM
9         dataset_flight
10    GROUP BY
11        1,2,3
12 )
13
14 SELECT
15     airline,
16     travel_month,
17     destination,
18     avg_price
19 FROM
20     MonthlyAverages
21 WHERE
22     rank = 1
23 ORDER BY
24     airline,
25     travel_month;
26

```

Completed

Time in queue: 127 ms

Run time: 760 ms

Data scanned: 1.26 MB

Results (38)

Copy

Download results

Q Search rows

< 1 >

⚙

#	airline	travel_month	destination	avg_price
1	Air Asia	03	Banglore	5049.444444444444
2	Air Asia	04	Banglore	4457.357142857143
3	Air Asia	05	Banglore	5125.357142857143
4	Air Asia	06	Banglore	5222.74358974359

PARQUET:

```

1 WITH MonthlyAverages AS (
2     SELECT airline,
3         split_part(date_of_journey, '/', 2) AS travel_month,
4         destination,
5         AVG(price) AS avg_price,
6     RANK() OVER (
7         PARTITION BY airline,
8             split_part(date_of_journey, '/', 2)
9         ORDER BY COUNT(*) DESC
10    ) AS rank
11 FROM flights_parquet
12 GROUP BY 1,
13     2,
14     3
15 )
16 SELECT airline,
17     travel_month,
18     destination,
19     avg_price
20 FROM MonthlyAverages
21 WHERE rank = 1
22 ORDER BY airline,
23     travel_month;

```

Completed

Time in queue: 105 ms

Run time: 672 ms

Data scanned: 47.56 KB

Results (38)

Copy

Download results

Search rows

< 1 > ⌕

#	airline	travel_month	destination	avg_price
1	Air Asia	03	Banglore	5049.444444444444
2	Air Asia	04	Banglore	4457.357142857143
3	Air Asia	05	Banglore	5125.357142857143
4	Air Asia	06	Banglore	5222.74358974359