IMPOPRTANT:

This report should only be read alongside the spark code, this PDF does not contain all the code, only the explanations and the outputs.

All Questions within a task in the python file should be ran separately, the single final task python file is mainly a guide (It does work but it is better to avoid trouble). (For some reason some questions cause conflict with each other).

All the work presented below is mine. As seen by the running time stamps in my screenshots, which (some) show my student's name: ec22954.

If you need to run a question yourself all the python files are named task(task_no.)-p(question no.) e.g. task2-p4.py will refer to the 4th question of task 2.

TASK 1:

Q1)

This question posed no challenge, I loaded the data frames with the data sets provided and printed the number of lines using the .count() method. However, utilising spark for the first time posed an exciting challenge. The code below reads: 22400728 entries.

```
2024-12-06 06:01:04,395 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 3: Stage finished
2024-12-06 06:01:04,396 INFO scheduler.DAGScheduler: Job 2 finished: count at NativeMethodAccessorImpl.java:0, took 22.639288 s
The number of entries is: 22400728
2024-12-06 06:01:04,664 INFO datasources.FileSourceStrategy: Pruning directories with:
2024-12-06 06:01:04,672 INFO datasources.FileSourceStrategy: Pushed Filters: IsNotNull(tpep_pickup_datetime),IsNotNull(trip_distance),GreaterThanOrEqual(tpep_al(tpep_pickup_datetime,2023-02-07)
2024-12-06 06:01:04,673 INFO datasources.FileSourceStrategy: Post-Scan Filters: isnotnull(tpep_pickup_datetime#40),isnotnull(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),isnotnull(trip_distance#43),(cast(fare_amount#47),
```

2024-12-06 06:01:04,674 INFO datasources.FileSourceStrategy: Output Data Schema: struct<tpep pickup datetime: string, trip distance: string, fare amount: string ... 1 more fields>

2024-12-06 06:01:04,394 INFO scheduler.DAGScheduler: ResultStage 3 (count at NativeMethodAccessorImpl.java:0) finished in 0.270 s 2024-12-06 06:01:04,395 INFO scheduler.DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job

e#43 as int) < 1),(tpep_pickup_datetime#40 >= 2023-02-01),(tpep_pickup_datetime#40 <= 2023-02-07)

Q2)

Applied the .filter() method to filter the data set to meet certain criteria. The code is pretty clear, the fare amount < 550, the trip_distance less than 1 and the data is the first week of august.

I then formatted the date using the date_format() function to convert to yyyy-MM-dd and used the GroupBy() function to count the pick-up dates. I believe this was the fastest and easiest way to approach this problem.

```
2024-12-06 06:01:43,115 INFO scheduler.DAGScheduler: ResultStage 13 (showString at NativeMethodAccessorImpl.java:0) finished in 0.480 s
2024-12-06 06:01:43,116 INFO scheduler DAGScheduler: Job 7 is finished. Cancelling potential speculative or zombie tasks for this job
2024-12-06 06:01:43,116 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 13: Stage finished
2024-12-06 06:01:43,116 INFO scheduler.DAGScheduler: Job 7 finished: showString at NativeMethodAccessorImpl.java:0, took 0.487279 s
2024-12-06 06:01:43,140 INFO codegen.CodeGenerator: Code generated in 15.530032 ms
 trip date trip count
2023-02-02
                   311
2023-02-06
                   344
 2023-02-01
                   275
2023-02-04
                   315
2023-02-03
                   318
2023-02-05
                   377
```

2024-12-06 06:01:43,115 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 13.0, whose tasks have all completed, from pool

```
2024-12-06 06:01:43,160 INFO server.AbstractConnector: Stopped Spark@466c94fd{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:01:43,161 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzfinaltaskl-p2-5-efb5a7939a8dd3d4-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:01:43,166 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:01:43,166 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
```

Q3)

I struggled quite a bit with this question at first. The join function was not working as intended at first, I believe it was because I was joining the wrong columns together so it was a human error which took me too long to fix.

One I have set the Location ID to their corresponding PULOcations and DOLOcations I dropped those and printed them.

However as seen by the code, I believe that I took the best approach, I joined the dataframes and manipulated the data efficiently.

The printSchematic() function was actually incredibly helpful throughout the project specially when it came to debugging and, understanding what exactly was in the data before playing around with it proved to be rather useful.

```
2024-12-06 05:07:44,962 INFO memory_MemoryStore: Block broadcast_5_piece0 stored as bytes in memory (estimated size 54.0 KiB, free 2003.5 MiB)
2024-12-06 05:07:44,963 INFO storage_BlockManagerInfo: Added broadcast_5_piece0 in memory on zzzzzzzzzzfinaltask1-p3-55db07939a5d6be6-driver-svc.data-science-ec22954.svc:7079 (size: 54.0 KiB, free: 2004.5 MiB)
2024-12-06 05:07:44,965 INFO spark.SparkContext: Created broadcast 5 from csv at NativeMethodAccessorImpl.java:0
2024-12-06 05:07:44,965 INFO execution.FileSourceScanExec: Planning scan with bin packing, max size: 134217728 bytes, open cost is considered as scanning 4194304 bytes.
```

```
-- Dropoff_service_zone: string (nullable = true)
2024-12-06 05:07:45,224 INFO server.AbstractConnector: Stopped Spark@4e99a17{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 05:07:45,226 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzfinaltaskl-p3-55db07939a5d6be6-driver-svc.data-science-ec22954.svc:4040 2024-12-06 05:07:45,230 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors 2024-12-06 05:07:45,231 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
```

Q4)

This question was quite easy, it was adding 2 new columns and using the month() method to extract only the month from the datetime value. The new columns created where route and month.

The use of concat_ws() method is certainly the best approach to tackle the use of route, whilst the method month() would pick-out the month from the data available. I believe there was no quicker way.

2024-12-06 06:07:03,800 INFO scheduler.DAGScheduler: Job 3 is finished. Cancelling potential speculative or zombie tasks for this job 2024-12-06 06:07:03,800 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 3: Stage finished 2024-12-06 06:07:03,800 INFO scheduler.DaGScheduler: Job 3 finished: shockString at NativeMethodAccessorImpl.java:0, took 4.995097 s 2024-12-06 06:07:03,805 INFO codegen.CodeGenerator: Code generated in 34.786616 ms |tpep_pickup_datetime|tpep_dropoff_datetime|passenger_count|trip_distance|payment_type|fare_amount|extra|mta_tax|tip_amount|tolls_amount|total_amount|congestion_surcharge|airport_fee|taxi_type |Pickup_Borough|Pickup_zone rough|Dropoff_zone |Dropoff_service_zone|route |Month| | 2023-01-01 00:32:10 | 2023-01-01 00:40:36 | 1.0 | Lenox Hill West | Yellow Zone | 2023-01-01 00:55:08 | 2023-01-01 01:01:27 | 1.0 | Upper East Side South | Yellow Zone 0.0 4.0 0.0 7.9 1.0 |0.5 16.9 2.5 |Upper East Side South|Yellow Zone |2023-01-01 00:25:04 |2023-01-01 00:37:49 |1.0 14.9 11.0 | 0.5 | 15.0 34.9 2.5 0.0 |Upper West Side North|Yellow Zone |2023-01-01 00:03:48 |2023-01-01 00:13:25 |0.0 |yellow_taxi|Queens 12.1 20.85 7.25 | 0.5 | 0.0 0.0 0.0 1.25 |LaGuardia Airport | Airports Queens to Queens |1 |1.0 |0.5 |3.28 |0.0 |yellow_taxi|Manhattan

nly showing top 5 rows

2024-12-06 06:07:03,877 INFO server.AbstractConnector: Stopped Spark#56b94faa{HTTP/1.1,{http/1.1}}{0.0.0.0:4040}
2024-12-06 06:07:03,879 INFO ui.SparkUI: Stopped Spark web UI at http://azzazzazzazzazzafaltaskl-p4-1-e7c083939a939f5b-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:07:03,885 INFO &8s.RubernetesClusterSchedulerBackend: Shutting down all executors

Q5)

Aggregated all the passengers with their associated tip_amount and calculated the average tip by dividing the total tip amount by the passenger count. A fairly easy question with only a few new functions such as agg().

2024-12-06 06:10:47,361 INFO scheduler.DAGScheduler: Job 4 is finished. Cancelling potential speculative or zombie tasks for this job 2024-12-06 06:10:47,361 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 6: Stage finished 2024-12-06 06:10:47,361 INFO scheduler.DAGScheduler: Job 4 finished: showString at NativeMethodAccessorImpl.java:0, took 0.212794 s 2024-12-06 06:10:47,387 INFO codegen.CodeGenerator: Code generated in 17.257272 ms

Month	route	passenger_count	total_tip_amount	average_tip_per_passenger	
+ 3 7 7 1 1 1 3 5	Manhattan to Queens Manhattan to Queens Manhattan to Unknown Queens to Bronx Bronx to Queens Unknown to EWR Staten Island to EWR	104937 89890 9902 6053 330 48	700702.1999999543 53321.17000000006 27825.890000000014 258.56 492.66999999999996 17.85	5.38488891133105 4.597041136626468 0.7835151515151515 10.26395833333333333333333333333333333333333	
6 4 12	Unknown to Manhattan Bronx to Bronx Manhattan to Manhattan	2665 2479 30	8531.2399999999994 1618.55 78.55	3.2012157598499043 0.6529043969342476 2.618333333333333	
+	+	, +	+	tt	

only showing top 10 rows

2024-12-06 06:10:47,406 INFO server.AbstractConnector: Stopped Spark@389afff9{HTTP/1.1,{http/1.1}}{0.0.0.0:4040}
2024-12-06 06:10:47,407 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzfinaltask1-p5-1-c963df939a961e6b-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:10:47,411 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:10:47,411 INFO k8s.KubernetesClusterSchedulerBackend\$KubernetesDriverEndpoint: Asking each executor to shut down

Q6 & 7)

I filtered the zero average tip by setting the average tip to 0. And top routes

Where ordered by ordering the average tip descending. Because it is OrderedBy() then the routes will also fall into order. This was my thinking process. OrderedBy() meant that by placing the average tip/p in descending order, it would automatically align the other rows with it. This meant that there is no need to order every column but that through ordering one, all others fall into place. This is also done later on and proves to be the most efficient way to organise data.

```
2024-12-06 06:15:18,884 INFO scheduler.DAGScheduler: ResultStage 24 (showString at NativeMethodAccessorImpl.java:0) finished in 1.295 s 2024-12-06 06:15:18,884 INFO scheduler.DAGScheduler: Job 14 is finished. Cancelling potential speculative or zombie tasks for this job 2024-12-06 06:15:18,885 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 24: Stage finished 2024-12-06 06:15:18,885 INFO scheduler.DAGScheduler: Job 14 finished: showString at NativeMethodAccessorImpl.java:0, took 54.722592 s 2024-12-06 06:15:18,910 INFO codegen.CodeGenerator: Code generated in 16.180877 ms
```

1	Month	route	passenger_count	total_tip_amount	average_tip_per_passenger
i	2	EWR to Brooklyn	1	35.76	35.76
İ	3	Staten Island to Unknown	4	122.66	30.665
İ	4	Bronx to EWR	1	30.51	30.51
İ	6	EWR to Brooklyn	2	56.93	28.465
ĺ	2	Bronx to EWR	2	55.26	27.63
İ	1	EWR to Queens	3	82.47	27.49
İ	1	Bronx to EWR	1	26.75	26.75
İ	1	EWR to Manhattan	2	46.3	23.15
İ	1	EWR to Staten Island	3	67.72	22.57333333333334
ĺ	5	Staten Island to Unknown	4	87.75	21.9375

only showing top 10 rows

2024-12-06 06:15:18,925 INFO server.AbstractConnector: Stopped Spark@1a78c89b{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:15:18,927 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzzinaltaskl-p6-1-556a65939a989581-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:15:18,932 INFO kBs.KubernetesClusterSchedulerBackend\$Shutting down all executors
2024-12-06 06:15:18,933 INFO kBs.KubernetesClusterSchedulerBackend\$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:15:18,942 WARN kBs.ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed (this is expected if the application is shutting down.)
2024-12-06 06:15:19,703 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!

TASK 2:

Q1)

Very simply I followed a straight forward approach similar to that of task 1. The printSchema() function was used to display the schema of the dataframe we are enquiring. It gives a clear view of the structure of the data, this includes column types, names, nulls...

This meant that I could understand any issues that could arise with the data that I want to manipulate before they came.

```
2024-12-06 06:21:24,210 INFO storage.BlockManagerInfo: Added broadcast 5 piece in memory on zzzzzzzzzzzzfinaltask2-p2-1-af2a96939aa0c7a7-driver-svc.data-science-ec22954.svc:7079 (size: 54.1 KiB, fre 2024-12-06 06:21:24,212 INFO spark.SparkContext: Created broadcast 5 from csv at NativeMethodAccessorImpl.java:0 2024-12-06 06:21:24,213 INFO execution.FileSourceScanExec: Planning scan with bin packing, max size: 93100933 bytes, open cost is considered as scanning 4194304 bytes.
    -- hash: string (nullable = true)
         nonce: string (nullable = true)
block_hash: string (nullable =
    -- block_number: string (nullable = true)
   -- transaction_index: string (nullable = true)
-- from_address: string (nullable = true)
-- to_address: string (nullable = true)
-- value: string (nullable = true)
-- gas_price: string (nullable = true)
-- input: string (nullable = true)
-- input: string (nullable = true)
-- block_timestamp: string (nullable = true)
-- max_fee_per_gas: string (nullable = true)
-- max_priority_fee_per_gas: string (nullable = true)
-- transaction_type: string (nullable = true)
    -- transaction index; string (nullable = true)
2024-12-06 06:21:24,424 INFO datasources.FileSourceStrategy: Pruning directories with: 2024-12-06 06:21:24,425 INFO datasources.FileSourceStrategy: Pushed Filters: 2024-12-06 06:21:24,425 INFO datasources.FileSourceStrategy: Post-Scan Filters:
2024-12-06 06:21:24,426 INFO datasources.FileSourceStrategy: Output Data Schema: struct<miner: string, size: string>
2024-12-06 06:21:24,470 INFO codegen.CodeGenerator: Code generated in 20.376268 ms
2024-12-06 06:21:24,504 INFO codegen.CodeGenerator: Code generated in 22.66602 ms
2024-12-06 06:21:24,510 INFO memory.MemoryStore: Block broadcast_6 stored as values in memory (estimated size 523.2 KiB, free 2002.9 MiB)
2024-12-06 06:21:24,533 INFO memory.MemoryStore: Block broadcast 6 piece0 stored as bytes in memory (estimated size 54.1 KiB, free 2002.9 MiB)
2024-12-06 06:21:24,535 INFO storage.BlockManagerInfo: Added broadcast 6 piece0 in memory on zzzzzzzzzzfinaltask2-p2-1-af2a96939aa0c7a7-driver-svc.data-science-ec22954.svc:7079 (size: 54.1 KiB, free 2004-12-06 06:21:24,536 INFO spark.SparkContext: Created broadcast 6 from showString at NativeMethodAccessorImpl.java:0
2024-12-06 06:21:24,541 INFO execution.FileSourceScanExec: Planning scan with bin packing, max size: 134217728 bytes, open cost is considered as scanning 4194304 bytes.
  2024-12-06 06:21:23,327 INFO memory MemoryStore: Block broadcast 2 piece0 stored as bytes in memory (estimated size 54.1 KiB, free 2003.5 MiB)
 2024-12-06 06:21:23,328 INFO storage.BlockManagerInfo: Added broadcast 2 from csv at NativeMethodAccessorImpl.java:0
2024-12-06 06:21:23,329 INFO spark.SparkContext: Created broadcast 2 from csv at NativeMethodAccessorImpl.java:0
2024-12-06 06:21:23,331 INFO execution.FileSourceScanExec: Planning scan with bin packing, max size: 134217728 bytes, open cost is considered as scanning 4194304 bytes.
    -- number: string (nullable = true)
     -- hash: string (nullable = true)
-- parent_hash: string (nullable = true)
-- nonce: string (nullable = true)
-- sha3_uncles: string (nullable = true)
     -- logs bloom: string (nullable = true)
-- transactions root: string (nullable = true)
-- state_root: string (nullable = true)
-- receipts_root: string (nullable = true)
    -- receipts root: string (nullable = true)
-- miner: string (nullable = true)
-- difficulty: string (nullable = true)
-- stze: string (nullable = true)
-- stze: string (nullable = true)
-- stra data: string (nullable = true)
-- gas limit: string (nullable = true)
-- gas used: string (nullable = true)
-- timestamp: string (nullable = true)
-- transaction count: string (nullable = true)
       -- transaction count: string (nullable = true)
    -- base_fee_per_gas: string (nullable = true)
  2024-12-06 06:21:23,520 INFO datasources.InMemoryFileIndex: It took 14 ms to list leaf files for 1 paths.
 2024-12-06 06:21:23,546 INFO datasources.InMemoryFileIndex: It took 8 ms to list leaf files for 1 paths.
 2024-12-06 06:21:23,704 INFO datasources.FileSourceStrategy: Pruning directories with:
2024-12-06 06:21:23,706 INFO datasources.FileSourceStrategy: Pushed Filters:
2024-12-06 06:21:23,707 INFO datasources.FileSourceStrategy: Post-Scan Filters: (length(trim(value#54, None)) > 0)
```

024-12-06 06:21:24,208 INFO memory.MemoryStore: Block broadcast 5 piece0 stored as bytes in memory (estimated size 54.1 KiB, free 2003.5 MiB)

Q2)

Here I used 2 important functions orderBy() which orders the columns of my dataframe with regards to a certain column. So in this case I orderedBy size so that I could get the top sizes mined and order the data frame by that standard. The other is desc() which simply orders them in descending order. (as asked). This is another example once again of being effective by re-arranging a single column and as that happens all other fall into place. Here:

```
0x52bc44d5378309e...
                             999
0x52bc44d5378309e...
                             999
0x77120710f4bf338...
                             999
0x1dcb8d1f0fcc8cb...
                             999
 0x52bc44d5378309e...
0x52bc44d5378309e...
                             999
0x52bc44d5378309e...
                             999
0x52bc44d5378309e...
                             999
0x52bc44d5378309e...
                             999
0x52bc44d5378309e...
only showing top 10 rows
2024-12-06 06:21:33,974 INFO server.AbstractConnector: Stopped Spark@69474243{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:21:33,976 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzfinaltask2-p2-1-af2a96939aa0c7a7-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:21:33,980 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:21:33,981 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
```

2024-12-06 06:21:33,896 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 2: Stage finished

2024-12-06 06:21:34,675 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!

2024-12-06 06:21:33,926 INFO codegen.CodeGenerator: Code generated in 19.13818 ms 2024-12-06 06:21:33,953 INFO codegen.CodeGenerator: Code generated in 18.43446 ms

miner total size

2024-12-06 06:21:33,897 INFO scheduler.DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 9.304849 s

As you see in the picture above, I made the mistake a couple of times that the show() function sets truncate to True by default. I later learnt to change this and my results become clearer. I learnt to be aware of this issue, not being able to visualise the data means poor testing or outright wrong output. Fixing this was essential and so I did:

2024-12-06 06:21:33,989 WARN k8s. Executor PodsWatch Snapshot Source: Kubernetes client has been closed (this is expected if the application is shutting down.)

```
2024-12-06 12:09:59,093 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 2: Stage finished
2024-12-06 12:09:59,093 INFO scheduler.DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 5.226119 s
2024-12-06 12:09:59,117 INFO codegen.CodeGenerator: Code generated in 13.668842 ms
2024-12-06 12:09:59,136 INFO codegen.CodeGenerator: Code generated in 13.657908 ms
lminer
                                            total size
0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 | 999
0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 | 999
0x77120710f4bf3383527dfb607a3b07b774276767 | 999
0x1dcb8d1f0fcc8cbc8c2d76528e877f915e299fbe | 999
0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 | 999
0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 999
 0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 999
0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 | 999
0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 | 999
0x52bc44d5378309ee2abf1539bf71de1b7d7be3b5 999
only showing top 10 rows
2024-12-06 12:09:59,200 INFO datasources.FileSourceStrategy: Pruning directories with:
2024-12-06 12:09:59,201 INFO datasources.FileSourceStrategy: Pushed Filters:
2024-12-06 12:09:59,201 INFO datasources.FileSourceStrategy: Post-Scan Filters:
```

Q3)

Here 2 very important function were used. date_format which successfully changed the format of the date to yyyy-MM-dd format. The other one which was imbedded within date_format, was from_unixtime which changed the format from a unix format to a regular date_format. Here is my output:

```
2024-12-06 06:23:20,832 INFO scheduler.DAGScheduler: ResultStage 2 (showString at NativeMethodAccessorImpl.java:0) finished in 0.413 s
2024-12-06 06:23:20,832 INFO scheduler DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job
2024-12-06 06:23:20,832 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 2: Stage finished
2024-12-06 06:23:20,832 INFO scheduler.DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 0.420840 s
2024-12-06 06:23:20,871 INFO codegen.CodeGenerator: Code generated in 23.537984 ms
 timestamp|formatted_date|
                1970-01-01
         0
1438269988
                2015-07-30
 1438270017
                2015-07-30
 1438270048
                2015-07-30
 1438270077
                2015-07-30
 1438270083
                2015-07-30
 1438270107
                2015-07-30
 1438270110
                2015-07-30
1438270112
                2015-07-30
1438270115
                2015-07-30
```

only showing top 10 rows

```
2024-12-06 06:23:20,899 INFO server.AbstractConnector: Stopped Spark@652c2443{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:23:20,902 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzzzfinaltask2-p3-1-3c9d8f939aa28738-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:23:20,908 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:23:20,909 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:23:20,901 WARN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed (this is expected if the application is shutting down.)
```

Q4)

This code aimed to perform an inner join between the 2 datasets provided by joining their similar hash columns. Here is the output for the number of entries: The join function has certainly been quite a challenge for me whilst working on this. I could not understand when an inner, outer, left or right join should take place. It was through trial and error that I understood the different outcomes each input gave.

```
2024-12-06 06:27:44,746 INFO scheduler.DAGScheduler: ResultStage 5 (count at NativeMethodAccessorImpl.java:0) finished in 0.080 s
2024-12-06 06:27:44,746 INFO scheduler.DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job
2024-12-06 06:27:44,746 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 5: Stage finished
2024-12-06 06:27:44,747 INFO scheduler.DAGScheduler: Job 2 finished: count at NativeMethodAccessorImpl.java:0, took 16.611026 s
The number of lines is: 504708
2024-12-06 06:27:44,764 INFO server.AbstractConnector: Stopped Spark@2dd3bc7d{HTTP/1.1, [http/1.1]}{0.0.0.0:4040}
2024-12-06 06:27:44,766 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzzzzfinaltask2-p4-1-d9b7ba939aa664fc-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:27:44,770 INFO k8s.KubernetesClusterSchedulerBackend: Shutzing down all executors
2024-12-06 06:27:44,771 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:27:44,778 WARN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed (this is expected if the application is shutting down.)
```

2024-12-06 06:27:44,744 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 5.0, whose tasks have all completed, from pool

Q5)

This section extracts the entries from only September by utilising the filter() function. By doing this a range between the following dates: col("formatted_date") >= "2015-09-01") & (col("formatted_date") < "2015-10-01" wav be used to only select the month of September. On top of this I grouped by formatted date and aggregated block_count and unique senders count number. Here is my output of the Datagraph (Incomplete):

```
2024-12-06 06:32:00,532 INFO scheduler.DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 52.056151 s 2024-12-06 06:32:00,554 INFO codegen.CodeGenerator: Code generated in 13.684788 ms 2024-12-06 06:32:00,570 INFO codegen.CodeGenerator: Code generated in 11.009585 ms
```

formatted_date	block_count	unique_senders_count_number
+	+	++
2015-09-01	1411	916
2015-09-02	1374	875
2015-09-03	1220	778
2015-09-04	1419	761
2015-09-05	1571	830
2015-09-06	1556	965
2015-09-07	1655	969
2015-09-08	1725	1003
2015-09-09	1732	980
2015-09-10	1685	989
2015-09-11	1870	1088
2015-09-12	1701	942
2015-09-13	1656	954
2015-09-14	1739	1013
2015-09-15	1673	906
2015-09-16	2050	920
2015-09-17	1883	982
2015-09-18	2035	891
2015-09-19	1786	840
2015-09-20	2206	904
+	+	++

only showing top 20 rows

```
2024-12-06 06:32:00,589 INFO server.AbstractConnector: Stopped Spark@718b6a5{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:32:00,591 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzzzzfinaltask2-p5-1-lfebc2939aa9bd6a-driver-svc.data-science-ec22954
2024-12-06 06:32:00,596 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:32:00,596 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:32:00,505 WARN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes Client has been closed (this is expected if the application is shutting
2024-12-06 06:32:01,397 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
```

Q6)

I primarily began by applying a filter to the formatted data so that it would only select the days of October, This makes it much easier to work with and gives opportunity for better quality outcomes / insights. I then performed a calculation by multiplying the gas column with the gas price column. This provided essential information regarding the total transaction fee of gas.

I lastly applied another filterwhich was used to refine the output data to only include rows where the transaction index was equal to 0. This made sure that the information being outputted stayed relevant to the desired criteria.

```
2024-12-06 06:36:11,065 INFO codegen.CodeGenerator: Code generated in 20.424756 ms
2024-12-06 06:36:11,093 INFO codegen.CodeGenerator: Code generated in 18.003883 ms
|formatted date|total transaction fee|
     2015-10-01 | 1.465063073190187...
     2015-10-02
                2.464981746013086E19
     2015-10-03 2.625289865558591E19
     2015-10-04 | 2.438632750193587E19
     2015-10-05 2.637662450505165...
     2015-10-06 2.418495276446883E19
     2015-10-07
                3.778585599630305E19
     2015-10-08 | 5.705912795259012E19
     2015-10-09 2.902720072304586...
     2015-10-10 2.201258740591100...
     2015-10-11 1.619052176793014...
     2015-10-12 2.168410861421186.
     2015-10-13 5.992645513202179E19
     2015-10-14 3.749759812719113E19
     2015-10-15 | 4.601426445537218E19
     2015-10-16
                 4.32299763950952E19
     2015-10-17 2.423809054826823..
     2015-10-18 | 2.876176820383356E19
     2015-10-19 2.647615515603579E19
     2015-10-20 3.402382771689090...
only showing top 20 rows
2024-12-06 06:36:11,119 INFO server.AbstractConnector: Stopped Spark@560c9579{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:36:11,121 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzfinaltask2-p6-1-d26f03939aadc05f-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:36:11,128 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:36:11,129 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:36:11,142 WARN k8s. Executor PodsWatch Snapshot Source: Kubernetes client has been closed (this is expected if the application is shutting down.)
```

2024-12-06 06:36:11,031 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 7: Stage finished

2024-12-06 06:36:12,228 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!

2024-12-06 06:36:11,031 INFO scheduler.DAGScheduler: Job 2 finished: showString at NativeMethodAccessorImpl.java:0, took 34.556880 s

TASK 3:

2024-12-06 06:36:12,253 INFO memory.MemoryStore: MemoryStore cleared 2024-12-06 06:36:12,254 INFO storage.BlockManager: BlockManager stopped

Q1)

Once Again, I printed the number of entries in the given Dataframe, not too difficult:

```
2024-12-06 06:41:01,305 INFO scheduler.DAGScheduler: ResultStage 3 (count at NativeMethodAccessorImpl.java:0) finished in 0.233 s 2024-12-06 06:41:01,306 INFO scheduler.DAGScheduler: Job 2 is finished. Cancelling potential speculative or zombie tasks for this job 2024-12-06 06:41:01,306 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 3: Stage finished 2024-12-06 06:41:01,307 INFO scheduler.DAGScheduler: Job 2 finished: count at NativeMethodAccessorImpl.java:0, took 1.703361 s The number of entries is: 466523 2024-12-06 06:41:01,537 INFO datasources.InMemoryFileIndex: It took 55 ms to list leaf files for 7 paths. 2024-12-06 06:41:01,601 INFO datasources.InMemoryFileIndex: It took 52 ms to list leaf files for 7 paths. 2024-12-06 06:41:01,712 INFO datasources.FileSourceStrategy: Pruning directories with: 2024-12-06 06:41:01,712 INFO datasources.FileSourceStrategy: Pushed Filters: 2024-12-06 06:41:01,712 INFO datasources.FileSourceStrategy: Post-Scan Filters: (length(trim(value#96, None)) > 0)
```

Q2)

Here I faced a challenge, since I was faced with issues when defining the StructTypes for the edges and vertecies. I did not quite understand what the purpose of this was. This meant that I was confused as to what I was doing. However. After going through the labs and enquiring online, the purpose for the use of these surprised me. They are higly applicable data processing techniques for modelling complex relationships between data point. It could prove to be useful for identifying interconnected nodes. Lab 7 however gave me the knowledge necessary for me to work on it and so I did. I firstly

defined the different StructFields within the StructType 'method' ensuring that my columns are the correct type. I then created 2 data frames and applied this logic to them. I used the show() function to display them. Here is the output:

2024-12-06 06:41:02,433 INFO scheduler TaskScheduler Impl: Killing all running tasks in stage 5: Stage finished

```
2024-12-06 06:41:02,434 INFO scheduler.DAGScheduler: Job 4 finished: showString at NativeMethodAccessorImpl.java:0, took 0.190642 s
2024-12-06 06:41:02,462 INFO codegen.CodeGenerator: Code generated in 20.847286 ms
|src|dst|
 |82 |196|
 236 229
only showing top 5 rows
2024-12-06 06:41:02,521 INFO datasources.FileSourceStrategy: Pruning directories with:
2024-12-06 06:41:02,522 INFO datasources.FileSourceStrategy: Pushed Filters:
2024-12-06 06:41:02,522 INFO datasources.FileSourceStrategy: Post-Scan Filters:
2024-12-06 06:41:02,522 INFO datasources.FileSourceStrategy: Output Data Schema: struct<LocationID: int, Borough: string, Zone: string, service_zone: string ... 2 mo
2024-12-06 06:41:02,822 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 6: Stage finished
2024-12-06 06:41:02.823 INFO scheduler DAGScheduler: Job 5 finished: showString at NativeMethodAccessorImpl.java:0. took 0.220655 s
2024-12-06 06:41:02,851 INFO codegen.CodeGenerator: Code generated in 18.560612 ms
|LocationID|Borough
                         Zone
                                                  service zone
11
            EWR
                          Newark Airport
                                                  EWR
12
                          Jamaica Bay
            Queens
                                                  Boro Zone
                          Allerton/Pelham Gardens Boro Zone
            Manhattan
                          Alphabet City
                                                  Yellow Zone
            Staten Island Arden Heights
15
                                                  Boro Zone
only showing top 5 rows
2024-12-06 06:41:02,867 INFO server.AbstractConnector: Stopped Spark@3b717763{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:41:02,869 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzfinaltask3-p2-1-ba021c939ab29f8f-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:41:02,876 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:41:02,877 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:41:02,887 WARN k8s. Executor PodsWatch Snapshot Source: Kubernetes client has been closed (this is expected if the application is shutting down.)
2024-12-06 06:41:03,712 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
```

Q3)

I then used the triplets function for analysing the relationship between the Pick up locations and the Drop off destinations. It was useful for displaying the vertices and edges that connect them. Gere was my output:

```
2024-12-06 06:53:01,643 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 5.0, whose tasks have all completed, from pool 2024-12-06 06:53:01,644 INFO scheduler.DASScheduler: ResultStage 5 (showString at NativeMethodAccessorImpl.java:0) finished in 0.268 s 2024-12-06 06:53:01,644 INFO scheduler.DASScheduler: Job 5 is finished. Cancelling potential speculative or zombie tasks for this job 2024-12-06 06:53:01,644 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 5: Stage finished 2024-12-06 06:53:01,644 INFO scheduler.DAGScheduler: Job 5 finished: showString at NativeMethodAccessorImpl.java:0, took 0.275627 s 2024-12-06 06:53:01,671 INFO codegen.CodeGenerator: Code generated in 19.089784 ms
   src
                                                                                                                                                                                                                                                                                                                                                                edge
                                                                                                                                                                                                                                                                                                                                                                                                                                     dst
                                                                                                                                                                                                                                                                                                                                                              [82, 196] [196, Queens, Rego Park, Boro Zone]
[7, 7] [7, Queens, Astoria, Boro Zone]
[7, 7] [7, Queens, Astoria, Boro Zone]
[17, 7] [7, Queens, Astoria, Boro Zone]
[16, 74] [74, Manhattan, Zast Harlem North, Boro Zone]
[236, 229] [229, Manhattan, Sutton Place/Turtle Bay North, Yellow Zone]
[75, 235] [235, Bronx, University Heights/Morris Heights, Boro Zone]
[260, 160] [160, Queens, Middle Village, Boro Zone]
[95, 264] [254, Unknown, NV, N/A]
[244, 41] [41, Manhattan, Central Harlem, Boro Zone]
[83, 7] [7, Queens, Astoria, Boro Zone]
[223, 223] [223, Queens, Steinway, Boro Zone]
[260, 260] [260, Queens, Woodside, Boro Zone]
[82, 193] [193, Queens, Queensbridge/Ravenswood, Boro Zone]
[146, 145] [145, Queens, Long Island City/Hunters Point, Boro Zone]
[130, 191] [191, Queens, Long Island City/Hunters Point, Boro Zone]
[195, 134] [134, Queens, Queens Village, Boro Zone]
[95, 134] [134, Queens, Queens Village, Boro Zone]
[195, 173] [173, Queens, North Corona, Boro Zone]
[196, 173] [173, Queens, North Corona, Boro Zone]
[197, 155] [155, Brooklyn, Marine Park/Mill Basin, Boro Zone]
   [82, Queens, Elmhurst, Boro Zone]
                                                                                                                                                                                                                                                                                                                                                                  [82, 196] [196, Queens, Rego Park, Boro Zone]
     [7, Queens, Astoria, Boro Zone]
[7, Queens, Astoria, Boro Zone]
 [[7, Queens, Astoria, Boro Zone]
[166, Manhattan, Morningside Heights, Boro Zone]
[236, Manhattan, Upper East Side North, Yellow Zone]
[75, Manhattan, East Harlem South, Boro Zone]
[260, Queens, Woodside, Boro Zone]
[95, Queens, Forest Hills, Boro Zone]
 [195, Queens, Forest Hills, Boro Zone]
[244, Manhattan, Washington Heights South, Boro Zone]
[83, Queens, Elmhurst/Maspeth, Boro Zone]
[250, Queens, Steinway, Boro Zone]
[82, Queens, Woodside, Boro Zone]
[82, Queens, Woodside, Boro Zone]
[82, Queens, Elmhurst, Boro Zone]
[146, Queens, Long Island City/Queens Plaza, Boro Zone]
[146, Queens, Long Island City/Queens Plaza, Boro Zone]
[130, Queens, Jamaica, Boro Zone]
[195, Queens, Forest Hills, Boro Zone]
[195, Queens, Forest Hills, Boro Zone]
[195, Queens, Forest Hills, Boro Zone]
[210, Brooklyn, Sheepshead Bay, Boro Zone]
```

```
2024-12-06 06:53:01,694 INFO server.AbstractConnector: Stopped Spark@3a04452e(HTTP[1.1,[http/1.1]]{0.0.0.0:4040}
2024-12-06 06:53:01,698 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzzzinaltask2-p3-4-3e4083939abda8aa-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:53:01,703 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:53:01,704 INFO k8s.KubernetesClusterSchedulerBackendSkubernetesEDriverEndpoint: Asking each executor to shut down
2024-12-06 06:53:01,718 WARN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes Client has been closed (this is expected if the application is shutting down.)
2024-12-06 06:53:02,504 INFO spark.MapOutputTrackerMasterEndpoint: Asking each executor to shut down
2024-12-06 06:53:02,504 INFO spark.MapOutputTrackerMasterEndpoint stopped!
2024-12-06 06:53:02,504 INFO spark.WashoryStore: MemoryStore cleared
```

This method was certainly insightful. I can more clearly see now after working with it myself the usefulness of creating connected data. It makes everything so much easier to understand and uncover much more information when diving into it. They help to uncover relationships and patterns that otherwise would be too obscure. And I certainly see the use of it too in within different employment sectors such as cyber security, and how connected data helps uncover suspeicious activity with greater ease.

Q4)

I then proceeded to count the connected vertices by using the find() function together with the count() function. After successfully counting these connected vertices I set up a dataframe. This dataframe includes both source and destination renamed in 2 columns as "id" and together with this their corresponding Boroughs and service_zones.

This creates a comprehensive view of the relationships between Borough and service_zone. This was incredible to work with and insightful. It proves to be a very efficient way of bringing out patterns from data.

```
2024-12-06 06:49:03,332 INFO scheduler.DAGScheduler: Job 8 is finished. Cancelling potential speculative or zombie tasks for this job
2024-12-06 06:49:03,332 INFO scheduler TaskScheduler Impl: Killing all running tasks in stage 9: Stage finished
2024-12-06 06:49:03,333 INFO scheduler.DAGScheduler: Job 8 finished: showString at NativeMethodAccessorImpl.java:0, took 0.287017 s
2024-12-06 06:49:03,363 INFO codegen.CodeGenerator: Code generated in 21.924467 ms
 id| id| Borough|service_zone
 82 | 196
            Queens
                      Boro Zone
            Oueens
                      Boro Zone
           Oueens
                      Boro Zone
     74 Manhattan
                      Boro Zone
236 229 Manhattan
                   Yellow Zone
260 | 160
           Queens
                      Boro Zone
 244 41 Manhattan
                      Boro Zone
 83
      7
            Queens
                      Boro Zone
223 223
            Queens
                      Boro Zone
260 260
            Queens
                      Boro Zone
```

only showing top 10 rows

```
2024-12-06 06:49:03,388 INFO server.AbstractConnector: Stopped Spark@3c96239e{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:49:03,391 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzzinaltask3-p4-1-0c46d9939ab9ec99-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:49:03,399 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:49:03,399 INFO k8s.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:49:03,391 WRN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed (this is expected if the application is shutting down.)
2024-12-06 06:49:04,232 INFO spark.MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint stopped!
```

Q5)

Here I used another very useful method in order to calculate the shortest path (no. of vertices) to reach a target destination. These target destination or target vertices are called landmark locations. I then made a table to illustrate theses result:

```
2024-12-06 06:55:48,818 INFO scheduler.DAGScheduler: ResultStage 64 (showString at NativeMethodAccessorImpl.java:0) finished in 0.160 s 2024-12-06 06:55:48,818 INFO scheduler.DAGScheduler: Job 14 is finished. Cancelling potential speculative or zombie tasks for this job
2024-12-06 06:55:48,818 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 64: Stage finished
2024-12-06 06:55:48,819 INFO scheduler.DAGScheduler: Job 14 finished: showString at NativeMethodAccessorImpl.java:0, took 0.167870 s
2024-12-06 06:55:48,838 INFO codegen.CodeGenerator: Code generated in 15.292641 ms
|id to 1|shortest distance
 147->1
 19->1
 39->1
            ĺ2
 71->1
 180->1
  130->1
  66->1
            2
 138->1
 171->1
             i2
 |170->1 |2
only showing top 10 rows
```

2024-12-06 06:55:48,856 INFO server.AbstractConnector: Stopped Spark@70fa8e23{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}

2024-12-06 06:55:48,859 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzinaltask3-p5-1-9e454e939abff184-driver-svc.data-science-ec22954.svc:4040 2024-12-06 06:55:48,869 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors 2024-12-06 06:55:48,870 INFO k8s.KubernetesClusterSchedulerBackend\$KubernetesDriverEndpoint: Asking each executor to shut down 2024-12-06 06:55:48,882 WARN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed (this is expected if the application is shutting down.)

Once again, it was quite instructive though challenging working with vertices and edges within data, to make things even as fundamental as finding the shortest path from a to b.

Q6)

This for a while forced me to read through pyspark documentation. However I was able to understand what it was that was being asked of me in this question. pageRank() is an important function within graphframes because they help to calculate the importance of vertices in a graph. The restProbability represents the probability of randomly jumping to another vertex instead of following the link, tol (tolerence) serves as a threshold for convergence and puts an end to the page ranking one the difference between iterations are smaller than the specified tol given. Here is the output for my program:

2024-12-06 06:58:19,129 INFO scheduler.TaskSchedulerImpl: Killing all running tasks in stage 2090: Stage finished

2024-12-06 06:58:19,161 INFO codegen.CodeGenerator: Code generated in 24.417212 ms

2024-12-06 06:58:19,130 INFO scheduler. DAGScheduler: Job 48 finished: showString at NativeMethodAccessorImpl.java:0, took 0.295073 s

2024-12-06 06:58:19,204 INFO server.AbstractConnector: Stopped Spark@6773f68e{HTTP/1.1,[http/1.1]}{0.0.0.0:4040}
2024-12-06 06:58:19,207 INFO ui.SparkUI: Stopped Spark web UI at http://zzzzzzzzzzzinaltask3-p6-1-78562b939ac20a58-driver-svc.data-science-ec22954.svc:4040
2024-12-06 06:58:19,215 INFO k8s.KubernetesClusterSchedulerBackend: Shutting down all executors
2024-12-06 06:58:19,216 INFO k8s.KubernetesClusterSchedulerBackend\$KubernetesDriverEndpoint: Asking each executor to shut down
2024-12-06 06:58:19,229 WARN k8s.ExecutorPodsWatchSnapshotSource: Kubernetes client has been closed (this is expected if the application is shutting down.)

The screenshots taken here do not have running time stamps since they are out of frame. I will include some at the end.

Q1) Here we apply the host and the port specifier to the session. This way we are able to set up a stable connection to the streaming data source. It is an essential step for stream data processing.

0024-12-06 07:43:42,515 INFO handler.ContextHandler: Started o.s. 024-12-06 07:43:42,517 INFO handler.ContextHandler: Started o.s. 024-12-06 07:43:42,517 INFO handler.ContextHandler: Started o.s. 02024-12-06 07:43:42,518 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0224-12-06 07:43:42,520 INFO handler.ContextHandler: Started o.s. 0	j.s.ServletContextHandl j.s.ServletContextHandl j.s.ServletContextHandl	ler@ef44eee(/SQL/json,null,AVAILABLE,@Spark} ler@lf7d05{/SQL/execution,null,AVAILABLE,@Spark} ler@55d2de4{/SQL/execution/json,null,AVAILABLE,@Spark}
Batch: 0		
t		
atch: 1		
	timestamp	
is working:,host,time,method,url,response,bytes 1,***.novo.dk,805465031,GET,/images/ksclogo-medium.gif,200,5866 0,***.novo.dk,805465029,GET,/ksc.html,200,7067 2,***.novo.dk,805465051,GET,/images/MOSAIC-logosmall.gif,200,363	2024-12-06 07:43:47.94 2024-12-06 07:43:49.94	46 16 18 17
Batch: 2		
logs tim	mestamp	
3,***.novo.dk,805465053,GET,/images/USA-logosmall.gif,200,234	24-12-06 07:43:50.95	

2024-12-06 07:43:42,489 INFO internal SharedState: Warehouse path is 'file:/opt/spark/work-dir/spark-warehouse'.

Q2)

In this question I implemented the use of the function Watermark().

Watermark is a mechanism that handles data that arrives late. Using the code written as an example with the following code:

logsDF = logsDF.withWatermark("timestamp", "3 seconds")

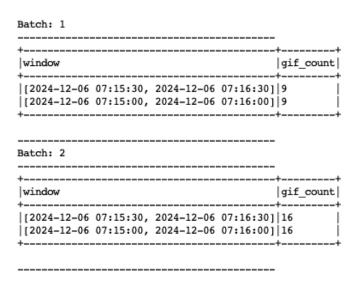
in this case what is taking place is that the logsDF dataframe with basis on the time stamp will allow late arriving data to bee processed but with a cap of 3 seconds late. This means that any data that arrives later than 3 seconds will not be taken in. It is a good way to manage late data whilst also limiting the waiting period on spark. There are many different uses for it which I have been reading about including: financial transactions, LoT streams and other applications that rely on event time.

Here some further components were added, this includes the watermark. This was explained in the section above, but to hammer the nail, the watermark is certainly an astonishing tool, by handling late incoming data, it maintains a consistency and accuracy that certainly make a difference in the results.

+	+-			+	+	+					+	+
logs +	t	imestamp		idx	hostna	ame t:	ime	method	resource		responsecode	byte +
20,***.novo.dk,807951864,GET,/images/NASA-logosmall.gif,2	00,786 2	024-12-06 07:	14:58.	451 20	***.nc	ovo.dk 8	07951864	GET	/images/NAS	A-logosmall.gi	£ 200	786
Batch: 20												
Batch: 20			+	+	+-		-+			+	+	
	timestan		+ idx	+ hostnai	+- ne t	time	-+ method	resourc	 ce	++ responsecode +	+ bytes	

Q4)

I have set up a window time and a slide time near the top of the python page, these are essential for configuring the behaviour of the window() function. The window function allows the grouping of streaming data into time-based windows. This way we can view and analyse patterns over specified intervals



Q5)

Once again the window() function is used, in this case it is used to group the streaming data based on hostnames and time intervals together with the total_bytes. In this occasion the total_bytes are calculated by the addition of the bytes column for each unique hostname instance and within a specified times period which was set to 30 seconds.

On top of this dropna() was utilised as a way to clean up the data so that it may resemble the examples provided. The original data consisted of some nulls. This improves quality and reliability of the input.

Batch: 4
+

+				+	++
window				hostname	total_bytes
+				+	++
[2024-12-06	07:28:00,	2024-12-06	07:29:00]	001.msy4.communique.net	44832.0
[2024-12-06	07:27:30,	2024-12-06	07:28:30]	***.novo.dk	247646.0
				001.msy4.communique.net	9630.0
[2024-12-06	07:27:00,	2024-12-06	07:28:00]	***.novo.dk	247646.0
[2024-12-06	07:27:30,	2024-12-06	07:28:30]	001.msy4.communique.net	54462.0
+				t	++

Q6)

The trigger functionality is utilised here, it plays a critical role when working with streaming data. It determines the frequency at which batch data is processed, and through doing this it massively increases it's efficiency by decreasing things like latency and throughput.

In the code we specify a trigger interval of 10 seconds. This means that the spark program will process batches of data at 10 second intervals. It means a consistent flow of incoming data which proves to be rather useful.

+	++
window	Correct count
+	++
007.thegap.com	19
001.msy4.communique.net	[8
***.novo.dk	29

Some running times stamps for question 4:

```
, task resources: Map(cpus -> name: cpus, amount: 1.0)
2024-12-06 07:34:37,784 INFO kbs.KubernetesClusterSchedulerBackend$KubernetesDriverEndpoint: Registered executor NettyRpcEndpointRef(spark-client://Executor) (10.134.94.30:38856) with ID 1
2024-12-06 07:34:37,889 INFO storage.BlockManagerMasterEndpoint: Registereing block manager 10.134.94.30:35473 with 2.1 GiB RAM, BlockManagerId(1, 10.134.94.30, 35473, None)
2024-12-06 07:34:38,122 INFO storage.BlockManagerMasterEndpoint: Registered executor NettyRpcEndpointRef(spark-client://Executor) (10.132.68.187:42038) with ID 2
2024-12-06 07:34:38,122 INFO kBs.KubernetesClusterSchedulerBackend$KubernetesdriverEndpoint: Registered executor NettyRpcEndpointRef(spark-client://Executor) (10.132.68.187:42038) with ID 2
2024-12-06 07:34:38,122 INFO kBs.KubernetesClusterSchedulerBackend: SchedulerBackend is ready for scheduling beginning after reached minRegisteredResourcesRatio: 0.8
2024-12-06 07:34:38,161 INFO storage.BlockManagerId(2, 10.132.68.187;39003 with 2.1 GiB RAM, BlockManagerId(2, 10.132.68.187, 39003, None)
2024-12-06 07:34:38,377 INFO internal.SharedState: Setting hive.metastore.warehouse.dir ('mill') to the value of spark.sql.warehouse.dir ('file:/opt/spark/work-dir/spark-warehouse')
2024-12-06 07:34:38,399 INFO internal.SharedState: Started o.s.j.s.ServletContextHandler@79f63aa8{/SQL,null,AVAILABLE, &Spark}}
2024-12-06 07:34:38,399 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@73665806(/SQL/sport,on,null,AVAILABLE, &Spark)
2024-12-06 07:34:38,399 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@73665806(/SQL/sport,on,null,AVAILABLE, &Spark)
2024-12-06 07:34:38,00 INFO handler.ContextHandler: Started o.s.j.s.ServletContextHandler@736066800d(/SQL/sport,on,null,AVAILABLE, &Spark)
 2024-12-06 07:34:38,400 INFO handler:ContextHandler: Started o.s.j.s.ServletContextHandler@17a3092{/SQL/execution/json,null,AVAILABLE,@Spark} 2024-12-06 07:34:38,402 INFO handler:ContextHandler: Started o.s.j.s.ServletContextHandler@2d178660{/static/sql,null,AVAILABLE,@Spark}
 |logs|timestamp|idx|hostname|time|method|resource|responsecode|bytes|
 Batch: 1
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         |method|resource
  logs
                                                                                                                                                                                                                                                                          |2024-12-06 07:34:42.921|Server
  Server
                                                                                                                                                                                                                                                                                                                                                                                                                null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 null
   is
                                                                                                                                                                                                                                                                          2024-12-06 07:34:42.921 is
                                                                                                                                                                                                                                                                                                                                                                                                                  null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            null
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 null
   working:,host,time,method,url,response,bytes | 2024-12-06 07:34:42.921 | working | 1,***.novo.dk,805465031,GET,/images/ksclogo-medium.gif,200,5866 | 2024-12-06 07:34:44.922 | 0,***.novo.dk,805465051,GET,/isc.html,200,7067 | 2024-12-06 07:34:43.921 | 0 | 2,****.novo.dk,805465051,GET,/images/MOSAIC-logosmall.gif,200,363 | 2024-12-06 07:34:45.923 | 2
                                                                                                                                                                                                                                                                                                                                                                                                                   host time meth

***.novo.dk 805465031 GET

***.novo.dk 805465029 GET

***.novo.dk 805465051 GET
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |bytes
|5866
|7067
|363
```

method url

|/images/ksclogo-medium.gii |200 |/ksc.html |200 |/images/MOSAIC-logosmall.gif|200

Batch: 3				
·				++
window				gif_count
[2024-12-06	07:35:30,	2024-12-06	07:36:30]	21
[2024-12-06	07:36:00,	2024-12-06	07:37:00]	9
Batch: 4				
+				++
window				gif_count
[2024-12-06	07:35:30,	2024-12-06	07:36:30]	30
[2024-12-06	07:36:00,	2024-12-06	07:37:00]	18

2024-12-06 07:36:43,040 ERROR v2.WriteToDataSourceV2Exec: Data source write support org.apache.spark.sql.execution.streaming.sources.MicroBatchWrite@24d53c07 is abort 2024-12-06 07:36:43,041 ERROR v2.WriteToDataSourceV2Exec: Data source write support org.apache.spark.sql.execution.streaming.sources.MicroBatchWrite@24d53c07 aborted

Batch: 8

+	+
window	Correct_count
4	

+	-++
007.thegap.com	34
001.msy4.communique.net	8
***.novo.dk	29
01-dynamic-c.wokingham.luna.ne	t 5

Batch: 9

+	++
window	Correct_count
+	++
007.thegap.com	34
001.msy4.communique.net	8
***.novo.dk	29
01-dynamic-c.wokingham.luna.net	15
+	++

2024-12-06 07:39:23,258 ERROR v2.WriteToDataSourceV2Exec: Data source write support org.apache.spark.sql.execution.streaming.sources.MicroBatchWrite@16ba981a is aborted.