# Index the Data Lake storage with Hyperspace

When loading data from Azure Data Lake Gen 2, searching in the data in one of the most resource consuming operations. Hyperspace introduces the ability for Apache Spark users to create indexes on their datasets, such as CSV, JSON, and Parquet, and use them for potential query and workload acceleration.
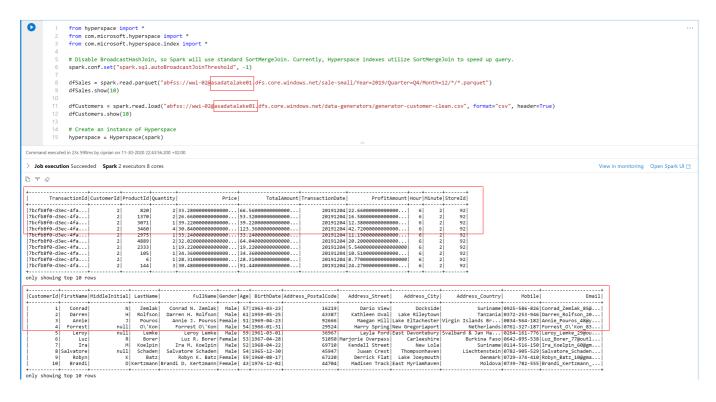
Hyperspace lets you create indexes on records scanned from persisted data files. After they're successfully created, an entry that corresponds to the index is added to the Hyperspace's metadata. This metadata is later used by Apache Spark's optimizer during query processing to find and use proper indexes. If the underlying data changes, you can refresh an existing index to capture that.

Also, Hyperspace allows users to compare their original plan versus the updated index-dependent plan before running their query.

Add a new cell to your notebook with the following code (remember to replace `<unique_suffix>` with the value you specified during the Synapse Analytics workspace deployment):

```python
from hyperspace import *
from com.microsoft.hyperspace import *
from com.microsoft.hyperspace.index import *

# Disable BroadcastHashJoin, so Spark will use standard SortMergeJoin. Currently,
Hyperspace indexes utilize SortMergeJoin to speed up query.
spark.conf.set("spark.sql.autoBroadcastJoinThreshold", -1)

dfSales = spark.read.parquet("abfss://wwi-
02@asagadatalake<unique_suffix>.dfs.core.windows.net/sale-
small/Year=2019/Quarter=Q4/Month=12/*/*.parquet")
dfSales.show(10)

dfCustomers = spark.read.load("abfss://wwi-
02@asagadatalake<unique_suffix>.dfs.core.windows.net/data-generators/generator-
customer-clean.csv", format="csv", header=True)
dfCustomers.show(10)

# Create an instance of Hyperspace
hyperspace = Hyperspace(spark)
```

Run the new cell. It will load the two DataFrames with data from the data lake and initalize Hyperspace.

```
1    from hyperspace import *
2    from com.microsoft.hyperspace import *
3    from com.microsoft.hyperspace.index import *
4
5    # Disable BroadcastHashJoin, so Spark will use standard SortMergeJoin. Currently, Hyperspace indexes utilize SortMergeJoin to speed up query.
6    spark.conf.set("spark.sql.autoBroadcastJoinThreshold", -1)
7
8    dfSales = spark.read.parquet("abfss://wwi-02@asadatalake01.dfs.core.windows.net/sale-small/Year=2019/Quarter=Q4/Month=12/*/*.parquet")
9    dfSales.show(10)
10
11   dfCustomers = spark.read.load("abfss://wwi-02@asadatalake01.dfs.core.windows.net/data-generators/generator-customer-clean.csv", format="csv", header=True)
12   dfCustomers.show(10)
13
14   # Create an instance of Hyperspace
15   hyperspace = Hyperspace(spark)
```

Command executed in 23s 590ms by ciprian on 11-30-2020 22:43:56.200 +02:00

> Job execution Succeeded  **Spark** 2 executors 8 cores                            View in monitoring   Open Spark UI

```
+---------------+----------+---------+--------+--------------------+--------------------+---------------+--------------------+----+------+-------+
|  TransactionId|CustomerId|ProductId|Quantity|               Price|         TotalAmount|TransactionDate|          ProfitAmount|Hour|Minute|StoreId|
+---------------+----------+---------+--------+--------------------+--------------------+---------------+--------------------+----+------+-------+
|7bcfb8f0-d3ec-4fa...|       2|      820|       2|33.28000000000000...|66.56000000000000...|       20191204|22.66000000000000...|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|     1370|       2|26.66000000000000...|53.32000000000000...|       20191204|16.58000000000000...|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|     3071|       1|39.22000000000000...|39.22000000000000...|       20191204|12.38000000000000...|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|     3460|       4|30.84000000000000...|123.36000000000000...|      20191204|42.72000000000000...|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|     2975|       1|33.24000000000000...|33.24000000000000...|       20191204|11.19000000000000...|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|     4889|       2|32.02000000000000...|64.04000000000000...|       20191204|20.20000000000000...|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|     2333|       1|19.22000000000000...|19.22000000000000...|       20191204|5.540000000000000000|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|      105|       1|34.36000000000000...|34.36000000000000...|       20191204|10.51000000000000...|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|        6|       1|28.31000000000000...|28.31000000000000...|       20191204|8.770000000000000000|   6|     2|     92|
|7bcfb8f0-d3ec-4fa...|       2|      144|       3|30.48000000000000...|91.44000000000000...|       20191204|24.27000000000000...|   6|     2|     92|
+---------------+----------+---------+--------+--------------------+--------------------+---------------+--------------------+----+------+-------+
only showing top 10 rows

+----------+---------+-------------+---------+--------------------+------+---+----------+--------------------+-------------+--------------+------------------+----------------+------------+--------------------+
|CustomerId|FirstName|MiddleInitial| LastName|            FullName|Gender|Age| BirthDate|Address_PostalCode|Address_Street|  Address_City|   Address_Country|          Mobile|       Email|
+----------+---------+-------------+---------+--------------------+------+---+----------+--------------------+-------------+--------------+------------------+----------------+------------+--------------------+
|         1|   Conrad|            N|   Zemlak|    Conrad N. Zemlak|  Male| 57|1963-03-23|             16219|   Dario View|      Dockside|          Suriname|0925-586-826|Conrad_Zemlak_85@...|
|         2|   Darren|            H|  Rolfson|   Darren H. Rolfson|  Male| 61|1959-05-25|             43387|Kathleen Oval| Lake Rileytown|          Tanzania|0372-253-946|Darren_Rolfson_20...|
|         3|    Annie|            J|   Pouros|    Annie J. Pouros|Female| 51|1969-04-23|             92666|  Maegan Hill|Lake Eltachester|Virgin Islands Br...|0834-564-182|Annie_Pouros_48@y...|
|         4|  Forrest|         null|  O\'Kon|    Forrest O\'Kon|  Male| 54|1966-01-31|             29524|  Harry Spring|New Gregoriaport|          Netherlands|0761-327-187|Forrest_O\'Kon_83...|
|         5|    Leroy|         null|    Lemke|         Leroy Lemke|  Male| 59|1961-03-01|             36967|   Layla Ford|East Davontebury|Svalbard & Jan Ma...|0264-161-776|Leroy_Lemke_29@ou...|
|         6|      Luz|            R|    Borer|     Luz R. Borer|Female| 53|1967-04-28|             51050|Marjorie Overpass|   Carleeshire|      Burkina Faso|0642-895-538|Luz_Borer_77@outl...|
|         7|      Ira|            M|  Koelpin|    Ira M. Koelpin|  Male| 52|1968-04-22|             69710|Kendall Street|     New Lola|          Suriname|0114-516-150|Ira_Koelpin_60@gm...|
|         8|Salvatore|         null|  Schaden|  Salvatore Schaden|  Male| 54|1965-12-30|             45947|  Juwan Crest| Thompsonhaven|      Liechtenstein|0782-905-529|Salvatore_Schaden...|
|         9|    Robyn|            K|     Batz|    Robyn K. Batz|Female| 59|1960-08-17|             67220| Derrick Flat|Lake Joeymouth|           Denmark|0729-374-410|Robyn_Batz_10@gma...|
|        10|   Brandi|            D|Kertzmann|Brandi D. Kertzmann|Female| 43|1976-12-02|             44704|Madisen Track|East Myriamhaven|           Moldova|0739-782-555|Brandi_Kertzmann_...|
+----------+---------+-------------+---------+--------------------+------+---+----------+--------------------+-------------+--------------+------------------+----------------+------------+--------------------+
only showing top 10 rows
```

Add another new cell to your notebook with the following code:

```
#create indexes: each one contains a name, a set of indexed columns and a set of
included columns
indexConfigSales = IndexConfig("indexSALES", ["CustomerId"], ["TotalAmount"])
indexConfigCustomers = IndexConfig("indexCUSTOMERS", ["CustomerId"], ["FullName"])

hyperspace.createIndex(dfSales, indexConfigSales)        # only create index
once
hyperspace.createIndex(dfCustomers, indexConfigCustomers)    # only create index
once
hyperspace.indexes().show()
```

Run the new cell. It will create two indexes and display their structure.

```
1    #create indexes: each one contains a name, a set of indexed columns and a set of included columns
2    indexConfigSales = IndexConfig("indexSALES", ["CustomerId"], ["TotalAmount"])
3    indexConfigCustomers = IndexConfig("indexCUSTOMERS", ["CustomerId"], ["FullName"])
4
5    hyperspace.createIndex(dfSales, indexConfigSales)        # only create index once
6    hyperspace.createIndex(dfCustomers, indexConfigCustomers)    # only create index once
7    hyperspace.indexes().show()
```

Command executed in 49s 512ms by ciprian on 11-30-2020 22:49:09.292 +02:00

> Job execution Succeeded  **Spark** 2 executors 8 cores                                                View in monitoring

```
+--------------+-------------+--------------+----------+--------------------+--------------------+------+
|          name|indexedColumns|includedColumns|numBuckets|              schema|       indexLocation| state|
+--------------+-------------+--------------+----------+--------------------+--------------------+------+
|indexCUSTOMERS| [CustomerId]|     [FullName]|       200|{"type":"struct",...|abfss://workspace...|ACTIVE|
|    indexSALES| [CustomerId]|  [TotalAmount]|       200|{"type":"struct",...|abfss://workspace...|ACTIVE|
+--------------+-------------+--------------+----------+--------------------+--------------------+------+
```

Add another new cell to your notebook with the following code:

```
df1 = dfSales.filter("""CustomerId = 203""").select("""TotalAmount""")
df1.show()
df1.explain(True)
```

Run the new cell. The output will show that the physical execution plan is not taking into account any of the indexes (performs a file scan on the original data file).



Now add another new cell to your notebook with the following code (notice the extra line at the beginning used to enable Hyperspace optimization in the Spark engine):

```
# Enable Hyperspace - Hyperspace optimization rules become visible to the Spark
optimizer and exploit existing Hyperspace indexes to optimize user queries
Hyperspace.enable(spark)
df1 = dfSales.filter("""CustomerId = 203""").select("""TotalAmount""")
df1.show()
df1.explain(True)
```

Run the new cell. The output will show that the physical execution plan is now using the index instead of the orginal data file.

Hyperspace provides an Explain API that allows you to compare the execution plans without indexes vs. with indexes. Add a new cell with the following code:

```
df1 = dfSales.filter("""CustomerId = 203""").select("""TotalAmount""")

spark.conf.set("spark.hyperspace.explain.displayMode", "html")
hyperspace.explain(df1, True, displayHTML)
```

Run the new cell. The output shows a comparison `Plan with indexes` vs. `Plan without indexes`. Observe how, in the first case the index file is used while in the second case the original data file is used.



Let's investigate now a more complex case, involving a join operation. Add a new cell with the following code:

```
eqJoin = dfSales.join(dfCustomers, dfSales.CustomerId ==
dfCustomers.CustomerId).select(dfSales.TotalAmount, dfCustomers.FullName)

hyperspace.explain(eqJoin, True, displayHTML)
```

Run the new cell. The output shows again a comparison `Plan with indexes` vs. `Plan without indexes`, where indexes are used in the first case and the original data files in the second.

In case you want to deactivate Hyperspace and cleanup the indexes, you can run the following code:

```
# Disable Hyperspace - Hyperspace rules no longer apply during query optimization.
# Disabling Hyperspace has no impact on created indexes because they remain intact
Hyperspace.disable(spark)

hyperspace.deleteIndex("indexSALES")
hyperspace.vacuumIndex("indexSALES")
hyperspace.deleteIndex("indexCUSTOMERS")
hyperspace.vacuumIndex("indexCUSTOMERS")
```