

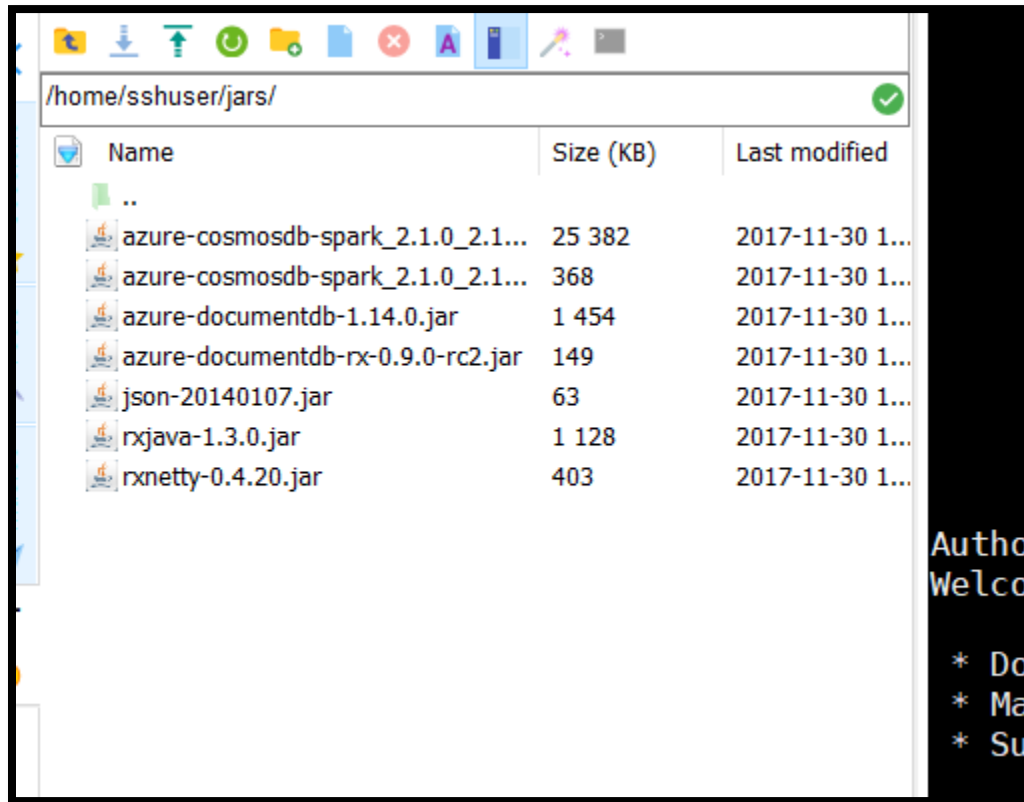
SPARK TO COSMOSDB

Setup Instruction

- 1) Uploading jar files to the cluster
- 2) Importing necessary jar files
- 3) Reading Data from CosmosDB
- 4) Writing Data to CosmosDB

Uploading jar files to the cluster

- 1) Download jar files from Maven repository
https://mvnrepository.com/artifact/com.microsoft.azure/azure-cosmosdb-spark_2.2.0_2.11/1.0.0
- 2) Upload jar files to the cluster



Reading Data from COSMOSDB

- 1) SSH to the cluster
- 2) Run Following command. This will initialize spark framework with reference of jar files mentioned

```
spark-shell --master yarn --jars /home/sshuser/jars/azure-cosmosdb-spark_2.1.0_2.11-1.0.0.jar,/home/sshuser/jars/azure-documentdb-1.14.0.jar,/home/sshuser/jars/rxjava-1.3.0.jar,/home/sshuser/jars/azure-documentdb-rx-0.9.0-rc2.jar,/home/sshuser/jars/json-20140107.jar
```

```
sshuser@hn0-sudhir:~$ spark-shell --master yarn --jars /home/sshuser/jars/azure-cosmosdb-
documentdb-1.14.0.jar,/home/sshuser/jars/rxjava-1.3.0.jar,/home/sshuser/jars/azure-docume
07.jar
SPARK_MAJOR_VERSION is set to 2, using Spark2
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel)
Spark context Web UI available at http://10.0.0.19:4040
Spark context available as 'sc' (master = yarn, app id = application_1512025053676_0004).
Spark session available as 'spark'.
Welcome to

  ____  _
 / ___|| | | |
| |___| |_| |
 \___ \|  _/
      |_|_|

 version 2.1.1.2.6.2.3-1

Using Scala version 2.11.8 (OpenJDK 64-Bit Server VM, Java 1.8.0_151)
Type in expressions to have them evaluated.
Type :help for more information.

scala>
```

- 3) Run following command to import reference and add configuration

```
import com.microsoft.azure.cosmosdb.spark.schema._
```

```
import com.microsoft.azure.cosmosdb.spark._
```

```
import com.microsoft.azure.cosmosdb.spark.config.Config
```

```
val ReadConfig = Config(Map("Endpoint" ->
```

```
"https://sudhirawdemocdb.documents.azure.com:443/",
```

```
"Masterkey" ->
```

```
"MstsCKBMXIh8zGfWx4YEhCUcTvBVor9yDYADpg7G393oqYHQX7erWFKE353cWsuqW3m9nZnxjyQ73lIpmeq==",
```

```
"Database" -> "startcosmos",
```

```
"preferredRegions" -> "Central India;",
```

```
"Collection" -> "assestdata",
```

```
"SamplingRatio" -> "1.0"))
```

```
scala> import com.microsoft.azure.cosmosdb.spark.schema._
import com.microsoft.azure.cosmosdb.spark.schema._

scala> import com.microsoft.azure.cosmosdb.spark._
import com.microsoft.azure.cosmosdb.spark._

scala> import com.microsoft.azure.cosmosdb.spark.config.Config
import com.microsoft.azure.cosmosdb.spark.config.Config

scala> val ReadConfig = Config(Map("Endpoint" -> "https://sudhirawdemocdb.documents.azure.com:443/",
| "Masterkey" -> "MstsCKBMXIh8zGfWx4YEhCUcTvBVor9yDYADpg7G393oqYHQX7erWFKE353cWsuqW3m9nZnxjyQ73lIpmeq==",
| "Database" -> "startcosmos",
| "preferredRegions" -> "Central India;",
| "Collection" -> "assestdata",
| "SamplingRatio" -> "1.0"))
ReadConfig: com.microsoft.azure.cosmosdb.spark.config.Config = com.microsoft.azure.cosmosdb.spark.config.ConfigBuilder$$anon$1@45ec2e46

scala>
```

- 4) Run following command to read data from the collection mentioned in configuration file

```
val coll = spark.sqlContext.read.cosmosDB(ReadConfig)
coll.createOrReplaceTempView("c")
val sqlDF = spark.sql("SELECT * FROM c ")
sqlDF.show(false)
```

```
scala> val coll = spark.sqlContext.read.cosmosDB(ReadConfig)
17/11/30 13:43:00 WARN CosmosDBConnection: CosmosDBConnection::Input preferred region list: Central India;
17/11/30 13:43:01 WARN ServiceJNIWrapper: 'Linux' with 'amd64' system is not compatible with native library. JNI not loaded.
coll: org.apache.spark.sql.DataFrame = [AlarmWord: string, AssetCode: string ... 36 more fields]

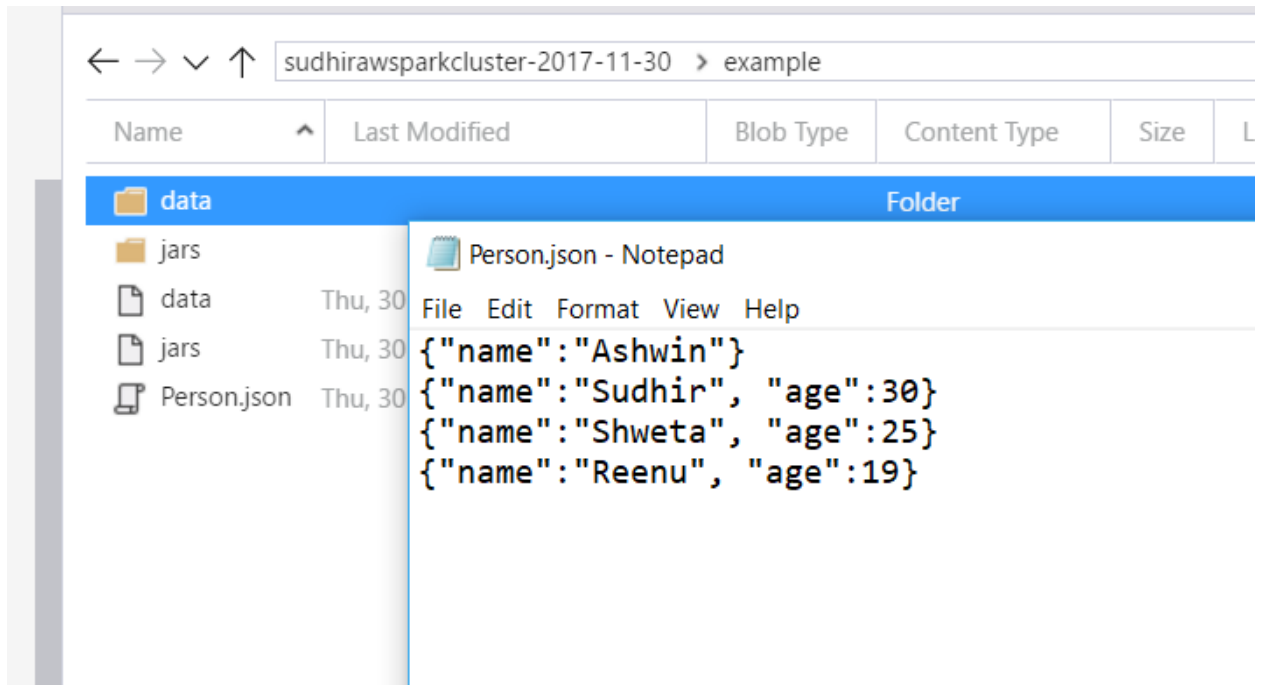
scala> coll.createOrReplaceTempView("c")

scala> val sqlDF = spark.sql("SELECT * FROM c ")
sqlDF: org.apache.spark.sql.DataFrame = [AlarmWord: string, AssetCode: string ... 36 more fields]

scala> sqlDF.show(false)
17/11/30 13:43:10 WARN Utils: Truncated the string representation of a plan since it was too large. This behavior can be adjusted by setting
'spark.debug.maxToStringFields' in SparkEnv.conf.
17/11/30 13:43:12 WARN CosmosDBConnection: CosmosDBConnection::Input preferred region list: Central India;
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|AlarmWord|AssetCode|AssetId|AssetSubCategory|CompanyCode|EventEnqueuedUtcTime|EventProcessedUtcTime|FuelType|GPSTimeStamp|Gatewa| | |
|HeightUnderHook|IMEI|InclinationX|InclinationY|InstanceOccuringDateTime|IoTHub|
|JibLength|Latitude|LoadValue|Longitude|MovementType|NooffFalls|PartitionId|RTTimeStamp|RatedLoad|ServerDateTime|SlewAngle|SoftwareVe|
sion|TrolleyPosition|WindSpeed|_attachments|_etag|_rid|_self|
|_ts|_age|_id|_name|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
|0|3450031|3450031|TC|1|2017-11-24T08:55:53.779Z|2017-11-24T08:58:11.0869452Z|A|35:34.0|RMT285|
42153|100|8.68326E+14|0|0|00:04.0|[d1,null,null,636471077750159279,2017-11-24T08:55:55.71|
Z,null]|5000|12.93342|280|77.73848|S|2|2|00:04.0|3000|24:03.0|149|V2.04
```

Writing Data to COSMOSDB

- 1) First upload sample json file with some dummy record to blob storage. These records will be inserted in COSMOSDB



- 2) Run following command

```
val persondf = spark.read.json("/example/Person.json")
```

```
scala> val persondf = spark.read.json("/example/Person.json")
persondf: org.apache.spark.sql.DataFrame = [age: bigint, name: string]
```

- 3) Run following command to point Azure cosmosdb (where data needs to be written)

```
val writeConfig = Config(Map("Endpoint" ->
  "https://sudhirawdemocdb.documents.azure.com:443/",
  "Masterkey" ->
  "MstsCKBMXIh8zGfWx4YEhCUcTvBVor9yDYADpg7G393oqYHQX7erWFKE353cWsu0qw3m9nZnx
  yjtvQ73lIpmeq==",
  "Database" -> "startcosmos",
  "PreferredRegions" -> "Central India;",
  "Collection" -> "assestdata",
  "WritingBatchSize" -> "100"))
```

```
scala> val writeConfig = Config(Map("Endpoint" -> "https://sudhirawdemocdb.documents.azure.com:443/",
  | "Masterkey" -> "MstsCKBMXIh8zGfWx4YEhCUcTvBVor9yDYADpg7G393oqYHQX7erWFKE353cWsu0qw3m9nZnx
  | "Database" -> "\nstartcosmos",
  | "PreferredRegions" -> "Central India;",
  | "Collection" -> "assestdata",
  | "WritingBatchSize" -> "100"))
writeConfig: com.microsoft.azure.cosmosdb.spark.config.Config = com.microsoft.azure.cosmosdb.spark.config.ConfigBuilder$$anon$1@2446dd1a
```

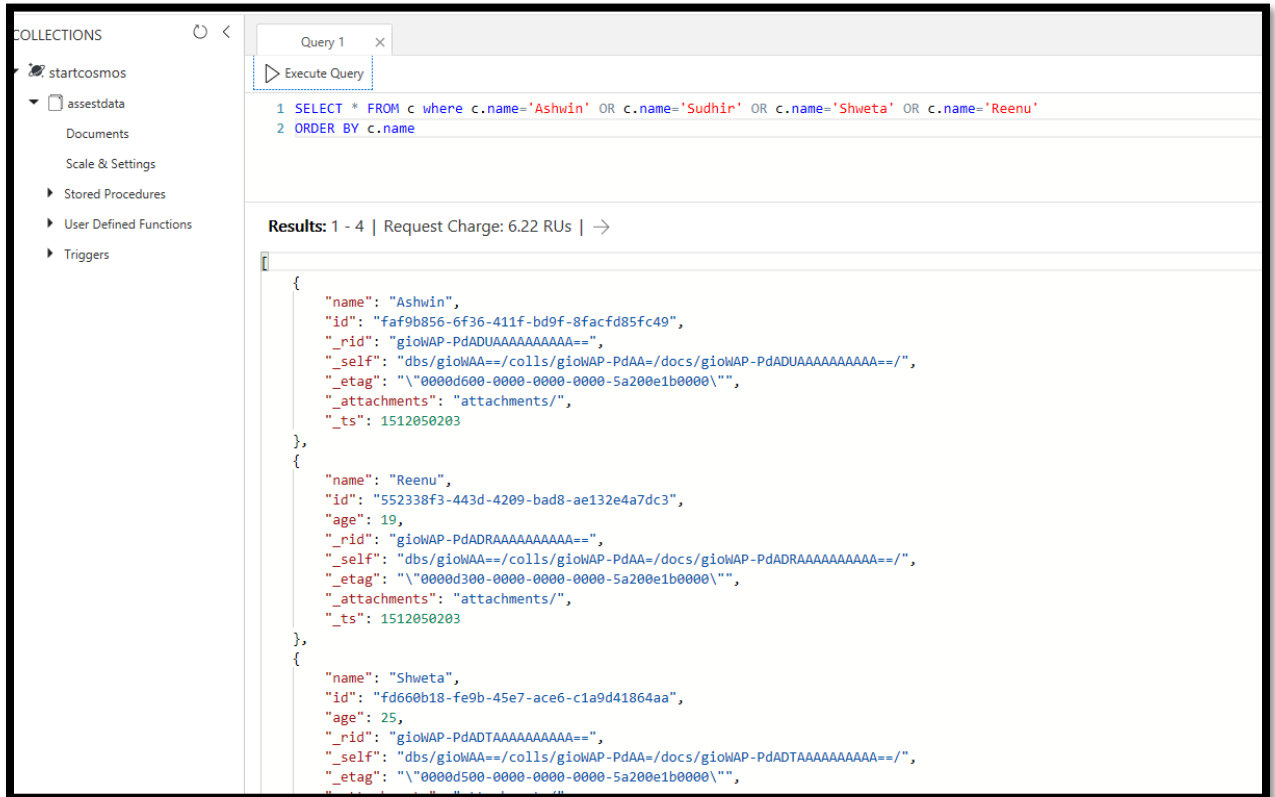
- 4) Run following command to Upsert dataframe

```
import org.apache.spark.sql.SaveMode
persondf.write.mode(SaveMode.Overwrite).cosmosDB(writeConfig)
```

```
scala> import org.apache.spark.sql.SaveMode
import org.apache.spark.sql.SaveMode

scala> persondf.write.mode(SaveMode.Overwrite).cosmosDB(writeConfig)
17/11/30 13:56:38 WARN CosmosDBConnection: CosmosDBConnection::Input preferred region list: Central India;
```

5) Let's read collection to make sure if data got inserted. From Azure portal



The screenshot shows the Azure Cosmos DB portal interface. On the left, the 'COLLECTIONS' sidebar lists 'startcosmos' and 'asstedata' (with sub-items: Documents, Scale & Settings, Stored Procedures, User Defined Functions, Triggers). The main area displays 'Query 1' with the following SQL query:

```
1 SELECT * FROM c where c.name='Ashwin' OR c.name='Sudhir' OR c.name='Shweta' OR c.name='Reenu'
2 ORDER BY c.name
```

Below the query, the results are shown as a JSON array of documents. The first document is for 'Ashwin' and the second is for 'Reenu'. The results section also indicates 'Results: 1 - 4 | Request Charge: 6.22 RUs | →'.

```
{
  "name": "Ashwin",
  "id": "faf9b856-6f36-411f-bd9f-8facfd85fc49",
  "_rid": "giowAP-PdADUAAAAAAAAA==",
  "_self": "dbs/giowAA=/colls/giowAP-PdAA=/docs/giowAP-PdADUAAAAAAAAA==/",
  "_etag": "\"0000d600-0000-0000-0000-5a200e1b0000\"",
  "_attachments": "attachments/",
  "_ts": 1512050203
},
{
  "name": "Reenu",
  "id": "552338f3-443d-4209-bad8-ae132e4a7dc3",
  "age": 19,
  "_rid": "giowAP-PdADRAAAAAAAAAA==",
  "_self": "dbs/giowAA=/colls/giowAP-PdAA=/docs/giowAP-PdADRAAAAAAAAAA==/",
  "_etag": "\"0000d300-0000-0000-0000-5a200e1b0000\"",
  "_attachments": "attachments/",
  "_ts": 1512050203
},
{
  "name": "Shweta",
  "id": "fd660b18-fe9b-45e7-ace6-c1a9d41864aa",
  "age": 25,
  "_rid": "giowAP-PdADTAAAAAAAAA==",
  "_self": "dbs/giowAA=/colls/giowAP-PdAA=/docs/giowAP-PdADTAAAAAAAAA==/",
  "_etag": "\"0000d500-0000-0000-0000-5a200e1b0000\"",
  "_attachments": "attachments/"
}
```