

#	Input	Output																		
1	<div><div>Data Description</div><div><div>➤ All the customer records are stored in the HDFS directory /user/spark/dataset/retail_db/customers-tab-delimited</div><div>➤ Data is in Text format</div><div>➤ Data is Tab delimited</div></div><div><div>Schema</div><table><tr><td>customer_id</td><td>int</td></tr><tr><td>customer_fname</td><td>string</td></tr><tr><td>customer_lname</td><td>string</td></tr><tr><td>customer_email</td><td>string</td></tr><tr><td>customer_password</td><td>string</td></tr><tr><td>customer_street</td><td>string</td></tr><tr><td>customer_city</td><td>string</td></tr><tr><td>customer_state</td><td>string</td></tr><tr><td>customer_zipcode</td><td>string</td></tr></table></div></div>	customer_id	int	customer_fname	string	customer_lname	string	customer_email	string	customer_password	string	customer_street	string	customer_city	string	customer_state	string	customer_zipcode	string	<div><div>Output Requirement</div><div><div>➤ Output all the customers who live in California</div><div>➤ Use text format for the output files</div><div>➤ Place the result data in /user/spark/dataset/result/scenario1/solution</div><div>➤ Result should only contain records that have state value as "CA"</div><div>➤ Output should only contain customer's full name Example: Robert Hudson</div></div><div></div></div>
customer_id	int																			
customer_fname	string																			
customer_lname	string																			
customer_email	string																			
customer_password	string																			
customer_street	string																			
customer_city	string																			
customer_state	string																			
customer_zipcode	string																			
2	<div><div>Data Description</div><div><div>➤ All the order records are stored in the HDFS directory /user/spark/dataset/retail_db/orders_parquet</div><div>➤ Data is in Parquet format</div></div><div><div>Output Requirement</div><div><div>➤ Output all the completed orders (where order_status value = "COMPLETE")</div><div>➤ Use JSON format for the output files</div><div>➤ Place the result data in HDFS directory /user/spark/dataset/result/scenario2/solution</div><div>➤ order_date should be in format yyyy-MM-dd</div><div>➤ Compress the output using gzip compression</div><div>➤ Output should only contain order_id, order_date,status</div></div><div></div></div></div>																			
3	<div><div>Data Description</div><div><div>➤ All the customer records are stored at /user/spark/dataset/retail_db/customers-tab-delimited</div><div>➤ Data is in text format</div><div>➤ Data is tab delimited</div></div><div><div>Schema</div><table><tr><td>customer_id</td><td>int</td></tr><tr><td>customer_fname</td><td>string</td></tr><tr><td>customer_lname</td><td>string</td></tr><tr><td>customer_email</td><td>string</td></tr><tr><td>customer_password</td><td>string</td></tr><tr><td>customer_street</td><td>string</td></tr><tr><td>customer_city</td><td>string</td></tr><tr><td>customer_state</td><td>string</td></tr><tr><td>customer_zipcode</td><td>String</td></tr></table></div></div>	customer_id	int	customer_fname	string	customer_lname	string	customer_email	string	customer_password	string	customer_street	string	customer_city	string	customer_state	string	customer_zipcode	String	<div><div>Output Requirement</div><div><div>➤ Output all the customers who live in "Caguas" city</div><div>➤ Place the result data in HDFS directory /user/spark/dataset/result/scenario3/solution</div><div>➤ Result should only contain records that have customer_city value as "Caguas"</div><div>➤ Compress the output using snappy compression</div><div>➤ Save the output using orc format</div></div><div></div></div>
customer_id	int																			
customer_fname	string																			
customer_lname	string																			
customer_email	string																			
customer_password	string																			
customer_street	string																			
customer_city	string																			
customer_state	string																			
customer_zipcode	String																			

4

Data Description

- All the categories records are stored at */user/spark/dataset/retail_db/categories*
- Data is in text format
- Data is comma separated

Schema

category_id	int
category_department_id	int
category_name	string

Output Requirement

- Convert data into tab delimited file
- Use text format for the output files
- Place the result data in HDFS directory */user/spark/dataset/result/scenario4/solution*
- Compress the output using lz4 compression

5

Data Description

- All the product records are stored at */user/spark/dataset/retail_db/products_avro*
- Data is in avro format
- Data is compressed with snappy compression

Output Requirement

- Output should only contain the products with price greater than 1000.0
- Use parquet format for the output files
- Place the result data in HDFS directory */user/spark/dataset/result/scenario5/solution*
- Compress the output using snappy compression

6

Data Description

- Get data from metastore table named "orders"
- Table is present in the database "default"

Output Requirement

- Fetch orders from Jan-2013 to Dec-2013
- Use parquet format for the output files
- Place the result data in HDFS directory */user/spark/dataset/result/scenario6/solution*
- Compress the output using Gzip compression

7

Data Description

- All the category records are stored at `/user/spark/dataset/retail_db/categories`
- Data is in text format
- Data is comma separated

Schema

category_id	int
category_department_id	int
category_name	string



Output Requirement

- Save all categories in metastore table `categories_replica` in default database
- Use no compression

8

Data Description

- All the category records are stored at `/user/spark/dataset/retail_db/categories`
- Data is in text format comma separated

Schema

category_id	int
category_department_id	int
category_name	string



Output Requirement

- create a metastore table named `'categories_parquet'`
- Table should only contain `category_id`, `category_name`
- Save all categories in metastore table `categories_parquet`

- Use parquet format for the output files

9

Data Description

- All the product records are stored at `/user/spark/dataset/retail_db/products_avro`
- Data is in avro format
- Data is compressed with snappy compression

NOTE:

- Use below command to start spark shell on Cloudera VM
`spark-shell --packages org.apache.spark:spark-avro_2.11:2.4.4` in CCA 175
- spark shell is already enabled with avro packages.

Output Requirement

- Output should contain columns `product_id`, `product_price`
- Save output as a JSON file
- Place the result data in HDFS directory `/user/spark/dataset/result/scenario9/solution`
- Use no compression

10

<div data-bbox="110 199 326 237">Data Description</div> <div data-bbox="110 254 604 308"><p>> All the category records are stored at <code>/user/spark/dataset/retail_db/categories</code></p><p>> Data is in text format comma separated</p></div> <div data-bbox="110 321 493 417"><div data-bbox="110 321 170 338">Schema</div><table><tr><td>category_id</td><td>int</td></tr><tr><td>category_department_id</td><td>int</td></tr><tr><td>category_name</td><td>string</td></tr></table></div> <div data-bbox="110 451 326 489">Output Requirement</div> <div data-bbox="110 497 501 583"><p>> Create a metastore table named 'categories_partitioned'</p><div data-bbox="110 527 464 548"></div><p>> Save all categories in metastore table categories_partitioned</p></div>	category_id	int	category_department_id	int	category_name	string	
category_id	int						
category_department_id	int						
category_name	string						