

Tasks to perform:

Week 1: Overview and basic configurations

Step 1: Choose a suitable cloud provider and set up a Spark shell environment

Step 2: Configure the necessary dependencies

Step 3: Execute basic Spark commands to make sure Spark is ready

Step 4: Use README.md for details, instructions, and commands

Week 2: Data ingestion

Step 1: Upload the entire data into Hive from CSV using cloud provider cluster setup (such as, EMR)

1. Log in to PuTTY with the username "hadoop"
2. Enter the command given below:

Command: hive

3. Create a database
4. Create a table with all the relevant details

Step 2: Create a bucket (for example: S3 and Azure Blob) and upload the csv file

Step 3: Load the data from the bucket into the Hive table

Step 4: Create a new directory in HDFS and copy the data from Hive into HDFS

Step 5: Check if the data has been successfully loaded in the HDFS path

Week 3: Data streaming

Step 1: Connect to Spark shell with all the dependencies (Hive, Hadoop, and HDFS).

1. Create Schema of the CSV files
2. Create a Spark session
 - Add Object Storage Service details as per the Cloud provider
 - Add all variables to your environment as they contain sensitive data

Step 2: Read the CSV file and convert the file to a data frame

Step 3: Convert "order_purchase_timestamp" to week and day using UDF

Step 4: Calculate the following data:

1. Total sales and order distribution per day and week for each city
2. Total sales and order distribution per day and week for each state
3. Average review score, average freight value, average order approval, and delivery time
4. The freight charges per city and total freight charges

Create a storage account

Project details

Select the subscription in which to create the new storage account. Choose a new or existing resource group to organize and manage your storage account together with other resources.

Subscription *

Simplilearn SS - 013

Resource group *

databricks-1393430

Create new

Instance details

Storage account name * ⓘ

retailanalytics

Region * ⓘ

(US) West US 2

Deploy to an Azure Extended Zone

Performance * ⓘ

☒ Standard: Recommended for most scenarios (general-purpose v2 account)

☐ Premium: Recommended for scenarios that require low latency.

Redundancy * ⓘ

Geo-redundant storage (GRS)

☒ Make read access to data available in the event of regional unavailability.

landing-data

Container

Search

Upload Change access level Refresh Delete Change tier Acquire lease Break lease View snapshots Create snapshot Give feedback

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Microsoft Entra user account)

Location: landing-data

Search blobs by prefix (case-sensitive)

Show deleted blobs

Add filter

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state	
<input type="checkbox"/> olist_public_dataset.csv	6/29/2024, 12:35:40 ...	Hot (Inferred)		Block blob	10.44 MiB	Available	***

Create a storage account ...

- Basics
- Advanced
- Networking
- Data protection
- Encryption
- Tags
- Review + create

[View automation template](#)

Basics

Subscription	Simplilearn SS - 013
Resource group	databricks-1393430
Location	West US 2
Storage account name	retailanalyticsstaging
Performance	Standard
Replication	Read-access geo-redundant storage (RA-GRS)

Advanced

Enable hierarchical namespace	Disabled
Enable SFTP	Disabled
Enable network file system v3	Disabled
Allow cross-tenant replication	Disabled
Access tier	Hot
Enable large file shares	Enabled

- Previous
- Next
- Create

Microsoft Azure

Search resources, services, and docs (G+)

Copilot

odf_user_1393430@sim...
SIMPLILEARN SHARED (SIMPLILEARN)

Home > retailanalyticsstaging_1719644813268 | Overview > retailanalyticsstaging

retailanalyticsstaging | Containers

Storage account

X

Search

+ Container

Change access level

Restore containers

Refresh

Delete

Give feedback

Search containers by prefix

Show deleted containers

Name	Last modified	Anonymous access level	Lease state	
<input type="checkbox"/> \$logs	6/29/2024, 12:38:05 PM	Private	Available	...
<input type="checkbox"/> staging-data	6/29/2024, 12:40:00 PM	Private	Available	...

[Home](#) >

Create Data Factory

[Basics](#) [Git configuration](#) [Networking](#) [Advanced](#) [Tags](#) [Review + create](#)

One-click to create data factory with sample pipeline and datasets. [Try it](#)

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription * ⓘ

Resource group * ⓘ [Create new](#)

Instance details

Name * ⓘ ✓

Region * ⓘ

Version * ⓘ

[Previous](#) [Next](#) [Review + create](#)

The screenshot displays the 'Review + create' step of the Azure Data Factory setup process. The interface is divided into several sections:

- Factory Resources:** A sidebar on the left showing a tree view of resources including Pipelines (Nullremover-pipeline), Datasets (DelimitedText1, olistoutput, olist), Data flows (dataflow1), and Power Query.
- Activities:** A central pane showing a 'Data flow' diagram with a 'NullRemover' activity. Below the diagram, there are tabs for 'General', 'Settings', 'Parameters', and 'User properties'. The 'General' tab is active, showing fields for Name (NullRemover), Description, Activity state (Activated), Timeout (0.12:00:00), Retry (0), and Retry interval (30).
- Properties:** A right-hand pane showing the 'General' tab for the 'Nullremover-pipeline' activity, with fields for Name and Description.
- Navigation:** At the top, there are buttons for 'Previous', 'Next', and 'Review + create'. Below these, there are tabs for 'Data Factory', 'Validate all', 'Publish all', and a 'Preview experience' toggle.

Factory Resources

Pipelines

- Nullremover-pipeline
- Change Data Capture (preview)

Datasets

- DelimitedText1
- olistoutput
- olist

Data flows

- dataflow1
- Power Query

Nullremover-pipeline

dataflow1

DelimitedText1

olistoutput

RawData

filtercol13

filtercol14

filtercol15

sink1

Columns: 15 total

Properties

General

Related (1)

Name

dataflow1

Description

Output stream name *

sink1

Description

Export data to DelimitedText2

Incoming stream *

filtercol15

Sink type *

Dataset

Inline

Cache

Dataset *

olistoutput

Skip line count

Microsoft Azure

Data Factory

sanir-datafactory

Search factory and documentation

odl_user_1393430@simplelearn.onmicrosoft.com

SimpleLearn Shared

All pipeline runs

Nullremover-pipeline - Activity runs

Activity runs

Pipeline run ID f58d939c-2ff5-4f63-b42f-ef9381d1c859

Activity status

Activity type

Run start

Duration

Integration runtime

User properties

Activity run ID

Log

Activity name	Activity status	Activity type	Run start	Duration	Integration runtime	User properties	Activity run ID	Log
NullRemover	Succeeded	Data flow	6/29/2024, 1:40:28 PM	3m 14s	AutoResolveIntegrationRu		c71bde46-a00a-49bc-dceb-aedc464777c	

Microsoft Azure

Search resources, services, and docs (G+/I)

Copilot

odl_user_1393430@sim...

SimpleLearn Shared (SIMPLE...

Home > retailanalyticsstaging | Containers >

staging-data

Container

Search

Upload

Change access level

Refresh

Delete

Change tier

Acquire lease

Break lease

View snapshots

Create snapshot

Give feedback

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Microsoft Entra user account)

Location: staging-data

Search blobs by prefix (case-sensitive)

Show deleted blobs

Add filter

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
_SUCCESS	6/29/2024, 1:43:39 PM	Hot (Inferred)		Block blob	0 B	Available
part-merged.csv	6/29/2024, 1:56:11 PM	Hot (Inferred)		Block blob	10.08 MiB	Available

Create an Azure Databricks workspace

- Basics
- Networking
- Encryption
- Security & compliance
- Tags
- Review + create

Project Details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *

Resource group *

Simplilearn SS - 013

databricks-1393430

Create new

Instance Details

Workspace name *

Region *

Pricing Tier *

Managed Resource Group name

databricksdeepi

West US 2

Standard (Apache Spark, Secure with Microsoft Entra ID)

Enter name for managed resource group

Home > deepi

deepi | Keys ☆ ...

Azure Cosmos DB account

Search Refresh Feedback

Overview
Activity log
Access control (IAM)
Tags
Diagnose and solve problems
Quick start
Notifications
Data Explorer

Settings

Features
Default consistency
Backup & Restore
Networking
CORS
Dedicated Gateway
Keys
Advisor Recommendations

URI
https://deepi.documents.azure.com:443/

Read-write Keys Read-only Keys

PRIMARY KEY
XRGXPoUdXjM3a4O2whj7LF3gRPWgy5a9SoYKIP5dCCt0wFRHvMFSWVeBvhHN586alr0CKztMPYnACDbo8u74g==
Last regenerated: 6/29/2024 (0 days ago). [Learn more](#)

SECONDARY KEY
Last regenerated: 6/29/2024 (0 days ago). [Learn more](#)

PRIMARY CONNECTION STRING
Last regenerated: 6/29/2024 (0 days ago). [Learn more](#)

SECONDARY CONNECTION STRING
Last regenerated: 6/29/2024 (0 days ago). [Learn more](#)

06:52 PM (1s) 1

```
# Define your storage account name and key
storage_account_name = "retailanalyticsstaging"
storage_account_key = "sjbGgHCFJbHuqOtTHSUBtGod+73z1DKhbX/6ac75CjBG+C37wTl2r0rKy2guYg8ECtm+FTqcrbim+ASTwQ0NfQ=="

# Define the container name and mount point
container_name = "staging-data"
mount_point = "/mnt/staging"

# List the files in the mounted directory to verify
display(dbutils.fs.ls(mount_point))
```

(3) Spark Jobs

	path	name	size	modificationTime
1	dbfs/mnt/staging/part-merged.c...	part-merged.c...	10574531	1719657496000

1 row | 1.27 seconds runtime

Refreshed 50 minutes ago

04:26 PM (9s) 2 Python

```
# Display the csv file as a spark dataframe
df = spark.read.csv("/mnt/staging/part-merged.csv")
```

(1) Spark Jobs

df: pyspark.sql.dataframe.DataFrame = [_c0: string, _c1: string ... 13 more fields]

04:27 PM (2s) 3

```
df.display(5)
```

(1) Spark Jobs

	_c0	_c1	_c2	_c3	_c4	_c5	_c6	_c7	_c8	_c9	_c10
1	1	delivered	79	17.8	1	Luzania	GO	728	50	201	2
2	2	delivered	119.9	27.16	1	Joinville	SC	892	50	511	3
3	3	delivered	519.99	41.69	1	Serra	ES	291	48	1156	2
4	4	delivered	28.5	17.92	1	RIO DE JANEIRO	RJ	222	21	207	2

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P databricksdeepi

New

- Workspace
- Recents
- Catalog
- Workflows
- Compute
- Data Engineering
- Job Runs
- Machine Learning
- Playground
- Experiments
- Features
- Models
- Serving
- Partner Connect

Catalog Explorer Send feedback

Type to filter

- hive_metastore
 - default
 - olistdatabase
 - olisttable
 - samples

olistdatabase >

olistdatabase.olisttable

Overview Sample Data Details History

Filter columns...

Column	Type	Comment
id	string	
order_status	string	
order_product_value	double	
order_freight_charge	double	
order_product_qty	int	
customer_city	string	
customer_state	string	
customer_zipcode_pre...	int	
product_name_len	int	
product_description_len	int	
product_photos_qty	int	

About this table

Data source format: CSV

Size: Unknown

Comment: Add comment

06:04 PM (1s) 8 Python

```
## Step15: Read the Hive table from the Hive database into a Spark DataFrame

df = spark.table("olistdatabase.olisttable")
df.display()
df.printSchema()
```

(1) Spark Jobs

df: pyspark.sql.dataframe.DataFrame = [id: string, order_status: string ... 13 more fields]

	id	order_status	order_product_value	order_freight_charge	order_product_qty	customer_city	customer_state	customer_zipcode_prefix
1	1	delivered	79	17.8	1	Luziania	GO	
2	2	delivered	119.9	27.16	1	Joinville	SC	
3	3	delivered	519.99	41.69	1	Serra	ES	
4	4	delivered	29.5	17.92	1	RIO DE JANEIRO	RJ	
5	5	delivered	26.77	23.11	1	Sao Paulo	SP	
6	6	delivered	419.9	23.02	1	Santa Adelia	SP	
7	7	delivered	65	16.21	1	Varginha	MG	
8	8	delivered	29.99	19.82	1	Sao Paulo	SP	
9	9	delivered	59.99	51.14	1	Carajas	PA	
10	10	delivered	56.99	16.13	1	Resende	RJ	
11	11	delivered	599	15.69	1	Guaratingueta	SP	
12	12	delivered	62	16.19	1	Tijucas	SC	
13	13	delivered	250	35.02	1	Mirassol	SP	
14	14	delivered	69.9	12.93	1	SAO PAULO	SP	

10,000+ rows | Truncated data due to row limit | 1.04 seconds runtime

Refreshed 1 hour ago

```
root
|-- id: string (nullable = true)
|-- order_status: string (nullable = true)
```

06:04 PM (<1s) 9 Python

```
from pyspark.sql import functions as F
from pyspark.sql.functions import col
from functools import reduce
```

06:04 PM (1s) 10

```
### Step16: Calculate the order metrics

df = df.withColumn("order_sales", col("order_product_value") * col("order_product_qty"))
df = df.withColumn("order_approval_time", col("order_approved_at") - col("order_placed_at"))
df = df.withColumn("order_delivery_time", col("order_delivered_at") - col("order_placed_at"))
df = df.withColumn("order_placed_date", F.to_date(col("order_placed_at")))
df = df.withColumn("order_placed_weekdate", F.trunc(col("order_placed_at"), "week"))
df.display();
```

(1) Spark Jobs

df: pyspark.sql.dataframe.DataFrame = [id: string, order_status: string ... 18 more fields]

	id	order_status	order_product_value	order_freight_charge	order_product_qty	customer_city	customer_state	customer_zipcode_prefix
1	1	delivered	79	17.8	1	Luziania	GO	
2	2	delivered	119.9	27.16	1	Joinville	SC	
3	3	delivered	519.99	41.69	1	Serra	ES	
4	4	delivered	29.5	17.92	1	RIO DE JANEIRO	RJ	
5	5	delivered	26.77	23.11	1	Sao Paulo	SP	
6	6	delivered	41.99	23.02	1	Santa Adelia	SP	
7	7	delivered	65	16.21	1	Varginha	MG	
8	8	delivered	29.99	19.82	1	Sao Paulo	SP	
9	9	delivered	59.99	51.14	1	Carajas	PA	
10	10	delivered	56.99	16.13	1	Resende	RJ	

06:04 PM (<1s) 11

```
# Step17: Compute the historical daily insights (We will need this values in further steps.)

grp_date = df.groupBy(col("order_placed_date"))
grp_city_date = df.groupBy(col("customer_city"), col("order_placed_date"))
grp_state_date = df.groupBy(col("customer_state"), col("order_placed_date"))
```

06:04 PM (<1s) 12

```
## Step18:

## Total sales
grp_date_total_sales = grp_date.agg(F.sum("order_sales").alias("total_sales")).orderBy(col("order_placed_date").asc())

## Total sales per city
grp_city_date_total_sales = grp_city_date.agg(F.sum("order_sales").alias("total_sales")).orderBy(col("customer_city").asc(), col("order_placed_date").asc())

## Total sales per state
grp_state_date_total_sales = grp_state_date.agg(F.sum("order_sales").alias("total_sales")).orderBy(col("customer_state").asc(), col("order_placed_date").asc())
```

grp_date_total_sales: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, total_sales: double]
grp_city_date_total_sales: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 1 more field]
grp_state_date_total_sales: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 1 more field]

06:04 PM (<1s) 13

```
## Step19:

#Total Freight Charge
#i. Total freight charge
#ii. Total freight charge per city
#iii. Total freight charge per state
## Total freight charge
grp_date_total_freight_charge = grp_date.agg(F.sum("order_freight_charge").alias("total_freight_charge")).orderBy(col("order_placed_date").asc())

## Total freight charge per city.
grp_city_date_total_freight_charge = grp_city_date.agg(F.sum("order_freight_charge").alias("total_freight_charge")).orderBy(col("customer_city").asc(), col("order_placed_date").asc())

## Total freight charge per state.
grp_state_date_total_freight_charge = grp_state_date.agg(F.sum("order_freight_charge").alias("total_freight_charge")).orderBy(col("customer_state").asc(), col("order_placed_date").asc())
```

grp_date_total_freight_charge: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, total_freight_charge: double]
grp_city_date_total_freight_charge: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 1 more field]
grp_state_date_total_freight_charge: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 1 more field]

06:04 PM (<1s) 14

```
## Step20:

# Total Order Count
# i. Total order count
# ii. Total order count per city
# iii. Total order count per state

## Total order count
grp_date_total_order_count = grp_date.agg(F.count("id").alias("total_order_count")).orderBy(col("order_placed_date").asc())

## Total order count per city
grp_city_date_total_order_count = grp_city_date.agg(F.count("id").alias("total_order_count")).orderBy(col("customer_city").asc(), col("order_placed_date").asc())

## Total order count per state
grp_state_date_total_order_count = grp_state_date.agg(F.count("id").alias("total_order_count")).orderBy(col("customer_state").asc(), col("order_placed_date").asc())
```

- grp_date_total_order_count: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, total_order_count: long]
- grp_city_date_total_order_count: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 1 more field]
- grp_state_date_total_order_count: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 1 more field]

06:04 PM (<1s) 15

```
# Step21: Average Freight charge

## Total average freight charge.
grp_date_avg_freight_charge = grp_date.agg(F.avg("order_freight_charge").alias("avg_freight_charge")).orderBy(col("order_placed_date").asc())

## Total order freight charge per city
grp_city_date_vg_freight_charge = grp_city_date.agg(F.avg("order_freight_charge").alias("avg_freight_charge")).orderBy(col("customer_city").asc(), col("order_placed_date").asc())

## Total order freight charge per state
grp_state_date_vg_freight_charge = grp_state_date.agg(F.avg("order_freight_charge").alias("avg_freight_charge")).orderBy(col("customer_state").asc(), col("order_placed_date").asc())
```

- grp_date_avg_freight_charge: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, avg_freight_charge: double]
- grp_city_date_vg_freight_charge: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 1 more field]
- grp_state_date_vg_freight_charge: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 1 more field]

06:04 PM (<1s) 16

```
# Step22: Average Review score

# Total average review charge.
grp_date_avg_review_score = grp_date.agg(F.avg("product_review_score").alias("avg_review_score")).orderBy(col("order_placed_date").asc())

## Total order average review per city
grp_city_date_avg_review_score = grp_city_date.agg(F.avg("product_review_score").alias("avg_review_score")).orderBy(col("customer_city").asc(), col("order_placed_date").asc())

## Total order average review per state
grp_state_date_avg_review_score = grp_state_date.agg(F.avg("product_review_score").alias("avg_review_score")).orderBy(col("customer_state").asc(), col("order_placed_date").asc())
```

- grp_date_avg_review_score: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, avg_review_score: double]
- grp_city_date_avg_review_score: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 1 more field]
- grp_state_date_avg_review_score: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 1 more field]

06:04 PM (<1s) 17

```
# Step23: Average Approval Time

grp_date_avg_approval_time = grp_date.agg(F.avg("order_approval_time").alias("avg_approval_time")).orderBy(col("order_placed_date").asc())

## Average Approval Time per city
grp_city_date_avg_approval_time = grp_city_date.agg(F.avg("order_approval_time").alias("avg_approval_time")).orderBy(col("customer_city").asc(), col("order_placed_date").asc())

## Average Approval Time per state
grp_state_date_avg_approval_time = grp_state_date.agg(F.avg("order_approval_time").alias("avg_approval_time")).orderBy(col("customer_state").asc(), col("order_placed_date").asc())
```

- grp_date_avg_approval_time: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, avg_approval_time: interval day to second]
- grp_city_date_avg_approval_time: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 1 more field]
- grp_state_date_avg_approval_time: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 1 more field]

06:04 PM (<1s) 18

```
# Step24: Average Delivery Time

grp_date_avg_delivery_time = grp_date.agg(F.avg("order_delivery_time").alias("avg_delivery_time")).orderBy(col("order_placed_date").asc())

## Average Approval Time per city
grp_city_date_avg_delivery_time = grp_city_date.agg(F.avg("order_delivery_time").alias("avg_delivery_time")).orderBy(col("customer_city").asc(), col("order_placed_date").asc())

## Average Approval Time per state
grp_state_date_avg_delivery_time = grp_state_date.agg(F.avg("order_delivery_time").alias("avg_delivery_time")).orderBy(col("customer_state").asc(), col("order_placed_date").asc())
```

- grp_date_avg_delivery_time: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, avg_delivery_time: interval day to second]
- grp_city_date_avg_delivery_time: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 1 more field]
- grp_state_date_avg_delivery_time: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 1 more field]

```
## ----- INSIGHTS TO BE STORED -----
# Historical daily insights into 3 tables

#I. Insights per period
#II. Insights per period per city
#II. Insights per period per state

## Per period
grp_date_insights = reduce(
    lambda x, y: x.join(y, on=["order_placed_date"], how="left"),
    [
        B"P_date_total_sales",
        B"P_date_total_freight_charge",
        B"P_date_total_order_count",
        B"P_date_avg_freight_charge",
        B"P_date_avg_review_score",
        B"P_date_avg_approval_time",
        B"P_date_avg_delivery_time",
    ]
)

## Per city
grp_city_date_insights = reduce(
    lambda x, y: x.join(y, on=["customer_city", "order_placed_date"], how="left"),
    [
        grp_date_insights,
        B"P_date_total_sales",
        B"P_date_total_freight_charge",
        B"P_date_total_order_count",
        B"P_date_avg_freight_charge",
        B"P_date_avg_review_score",
        B"P_date_avg_approval_time",
        B"P_date_avg_delivery_time",
    ]
)
```

```
## Per State
grp_state_date_insights = reduce(
    lambda x, y: x.join(y, on=["customer_state", "order_placed_date"], how="left"),
    [
        B"P_state_date_total_sales",
        B"P_state_date_total_freight_charge",
        B"P_state_date_total_order_count",
        B"P_state_date_avg_freight_charge",
        B"P_state_date_avg_review_score",
        B"P_state_date_avg_approval_time",
        B"P_state_date_avg_delivery_time",
    ]
)
```

▶ `grp_date_insights`: pyspark.sql.dataframe.DataFrame = [order_placed_date: date, total_sales: double ... 6 more fields]
▶ `grp_city_date_insights`: pyspark.sql.dataframe.DataFrame = [customer_city: string, order_placed_date: date ... 7 more fields]
▶ `grp_state_date_insights`: pyspark.sql.dataframe.DataFrame = [customer_state: string, order_placed_date: date ... 7 more fields]

```
display(grp_city_date_insights)
```

▶ (14) Spark Jobs

	customer_city	order_placed_date	1.2 total_sales	1.2 total_freight_charge	1.2 total_order_count	1.2 avg_freight_charge	1.2 avg_review_score	1.2 avg_approval_time
1	Sao Paulo	2017-11-27	7866.610000000001	808.9899999999999	51	15.86254901960784	4	["seconds":2
2	Sao Goncalo	2017-01-30	19.99	16.32	1	16.32	5	["seconds":8

Week 4: Data analysis and visualization

Step 1: Write the results into HDFS

Step 2: Save the final dataset into object storage service per the cloud platform

Step 3: Create a DB cluster that is also a NoSQL using the relevant service on the cloud platform

Step 4: Save insights in the NoSQL DB mentioned in the previous step

Remove “Enable soft delete for blobs” and “Enable soft delete for containers”

Home > retailanalytics

retailanalytics | Data protection

Storage account

Search

- Overview
- Activity log
- Tags
- Diagnose and solve problems
- Access Control (IAM)
- Data migration
- Events
- Storage browser
- Storage Mover
- Data storage
- Security + networking
- Data management
 - Storage tasks (preview)
 - Redundancy
 - Data protection**
 - Object replication
 - Blob inventory

Data protection provides options for recovering your data when it is erroneously modified or deleted.

Recovery

- ☐ Enable Azure Backup for blobs
- ☐ Enable point-in-time restore for containers
- ☐ Enable soft delete for blobs
- ☐ Enable soft delete for containers
- ☐ Enable permanent delete for soft deleted items

Tracking

- ☐ Enable versioning for blobs
- ☐ Enable blob change feed

Access control

- ☐ Enable version-level immutability support

Save Discard

```
# Step27: Write the insights as a CSV file into the file system [DBFS]

# write pyspark dataframe as csv file into file system (dbfs -> Databricks file system associated with every dbricks workspace.)
grp_date_insights.write.csv("dbfs://FileStore/shared_uploads/odl_user_1393560@simplelearnss.onmicrosoft.com/grpdateinsights", mode="overwrite", header=True)
grp_city_date_insights.write.csv("dbfs://FileStore/tables/grpcitydateinsights", mode="overwrite", header=True)
grp_state_date_insights.write.csv("dbfs://FileStore/tables/grpstatedateinsights", mode="overwrite", header=True)

(42) Spark Jobs

# Step 28 : ADLS Blob Object storage
# write pyspark dataframe as csv file into storage account (Blob Storage Azure/AWS protocol)
# Set the Azure storage account access key
spark.conf.set(
    "fs.azure.account.key.retailanalyticsstaging.dfs.core.windows.net",
    "sJbGgHCF7bhuqDtLTHSUBtG0d+73zLDKbX/6ac75Cj9G+C37wTl2r0rky2guVg8ECTm+FTqcrbm+ASTwQ8NFQ=="
)

grp_date_insights.write.csv(
    "abfss://staging-data@retailanalyticsstaging.dfs.core.windows.net/insights/grpdateinsights.csv",
    mode="overwrite", header=True
)

grp_city_date_insights.write.csv(
    "abfss://staging-data@retailanalyticsstaging.dfs.core.windows.net/insights/grpcitydateinsights.csv",
    mode="overwrite", header=True
)

grp_state_date_insights.write.csv(
    "abfss://staging-data@retailanalyticsstaging.dfs.core.windows.net/insights/grpstatedateinsights.csv",
    mode="overwrite", header=True
)

(42) Spark Jobs
```

staging-data

Search

Upload

Change access level

Refresh

Delete

Change tier

Acquire lease

Break lease

View snapshots

Create snapshot

Give feedback

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Microsoft Entra user account)

Location: staging-data

Search blobs by prefix (case-sensitive)

Show deleted blobs

Add filter

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
grpdatainsights						-
insights						-
grpdatainsights	6/29/2024, 6:58:59 PM	Hot (Inferred)		Block blob	0 B	Available
insights	6/29/2024, 7:01:28 PM	Hot (Inferred)		Block blob	0 B	Available
part-merged.csv	6/29/2024, 4:08:16 PM	Hot (Inferred)		Block blob	10.08 MiB	Available

staging-data

Search

Upload

Change access level

Refresh

Delete

Change tier

Acquire lease

Break lease

View snapshots

Create snapshot

Give feedback

Overview

Diagnose and solve problems

Access Control (IAM)

Settings

Shared access tokens

Access policy

Properties

Metadata

Authentication method: Access key (Switch to Microsoft Entra user account)

Location: staging-data / insights

Search blobs by prefix (case-sensitive)

Show deleted blobs

Add filter

Name	Modified	Access tier	Archive status	Blob type	Size	Lease state
[.]						-
grpcitydateinsights.csv						-
grpdatainsights.csv						-
grpstatedateinsights.csv						-
grpcitydateinsights.csv	6/29/2024, 7:01:30 PM	Hot (Inferred)		Block blob	0 B	Available
grpdatainsights.csv	6/29/2024, 7:01:28 PM	Hot (Inferred)		Block blob	0 B	Available
grpstatedateinsights.csv	6/29/2024, 7:01:32 PM	Hot (Inferred)		Block blob	0 B	Available

Microsoft Azure

Search resources, services, and docs (G+)

Copilot

Home > Azure Cosmos DB >

Create Azure Cosmos DB Account - Azure Cosmos DB for NoSQL

Project Details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription *

Simplilearn SS - 014

Resource Group *

databricks-1393560

Create new

Instance Details

Account Name *

deepi

Configure availability zone settings for your account. You cannot change these settings once the account is created.

Availability Zones ⓘ

Enable

Disable

Location * ⓘ

(US) West US 3

Available locations are determined by your subscription's access and availability zone support (if that is enabled). If you don't see or cannot select your desired location access.
[Click here for more details on how to create a region access request](#)

Capacity mode ⓘ

Provisioned throughput

Serverless

[Learn more about capacity mode](#)

Review + create

Previous

Next: Global distribution

Create Azure Cosmos DB Account - Azure Cosmos DB for NoSQL ...

Validation Success

- Basics
- Global distribution
- Networking
- Backup Policy
- Encryption
- Tags
- Review + create

Creation Time

Estimated Account Creation Time (in minutes)

3

The estimated creation time is calculated based on the location you have selected

Basics

Subscription	Simplilearn SS - 014
Resource Group	databricks-1393560
Location	West US 3
Account Name	(new) deepi
API	Azure Cosmos DB for NoSQL
Capacity mode	Serverless
Availability Zones	Disable

Backup Policy

Create

Previous

Next

Download a template for automation

Deployment

Search

Delete

Cancel

Redeploy

Download

Refresh

Overview

Inputs

Outputs

Template

✔ Your deployment is complete

Deployment name : Microsoft.Azure.CosmosDB-20240629115632

Subscription : Simplilearn SS - 013

Resource group : databricks-1393430

Start time : 6/29/2024, 11:56:44 AM

Correlation ID : b3d5a152-4054-4ddf-b897-e22ddb4e7707

> Deployment details

< Next steps

Go to resource

Give feedback

Tell us about your experience with deployment

1. Navigate to your Databricks workspace.
2. Go to the "Clusters" section and select the cluster where you want to install the connector.
3. In the cluster configuration, find the "Libraries" section and click on "Install New".
4. Choose "Maven" as the source for the library.
5. In the search box, enter the Maven coordinates for the Azure Cosmos DB Spark connector (com.azure.cosmos.spark:azure-cosmos-spark_3-4:4.0)

Compute >

Samir's Cluster

Configuration Notebooks (0) **Libraries** Event log Spark UI Driver logs Metrics Apps Spark compute UI - Master ▾

Filter libraries

Uninstall Install new

<input type="checkbox"/>	Status	Name	Type	Source
<input type="checkbox"/>		com.azure.cosmos.spark:azure-cosmos-spark_3-1_2-12:4.11.2	Maven	-
<input type="checkbox"/>		com.azure.cosmos.spark:azure-cosmos-spark_3-4:4.0	Maven	-

1

```

10:42 PM (<1s) 22
from pyspark.sql.functions import col, expr

# Assuming avg_approval_time is of type INTERVAL DAY TO SECOND, convert it to seconds
grp_date_insights = grp_date_insights.withColumn(
    "avg_approval_seconds",
    expr("CAST(avg_approval_time AS LONG)")
)

grp_date_insights = grp_date_insights.withColumn(
    "avg_delivery_seconds",
    expr("CAST(avg_delivery_time AS LONG)")
)

grp_date_insights: pyspark.sql.dataframe.DataFrame = [order_placed_weekdate: date, total_sales: double ... 9 more fields]

10:46 PM (2s) 23
# Drop unwanted columns
grp_date_insights = grp_date_insights.drop("avg_approval_time", "avg_delivery_time", "extracted_seconds")
display(grp_date_insights)

(14) Spark Jobs
grp_date_insights: pyspark.sql.dataframe.DataFrame = [order_placed_weekdate: date, total_sales: double ... 6 more fields]

Table +

```

```

07:40 PM (<1s) 23

## Step28: Write the insights to NoSQL

# template
"""
config = {
    "spark.cosmos.accountEndpoint": "<cosmos-account-endpoint>",
    "spark.cosmos.accountKey": "<cosmos-account-key>",
    "spark.cosmos.database": "<database-name>",
    "spark.cosmos.container": "<container-name>",
}
"""

# connect to cosmos account using read-write key
config = {
    "spark.cosmos.accountEndpoint": "https://deepi.documents.azure.com:443/",
    "spark.cosmos.accountKey":
    "XRGPoUjXjI3a402whj7LF3gRPHgy5a9SoYKJP5dCct8xFRhVhF5WVeBvHhN586a1r0CKztHPTyNACDboBu74g==",
    "spark.cosmos.database": "ToDoList",
    "spark.cosmos.container": "Items",
}

```


▶ ✓ 10:52 PM (5s)

25

```
from pyspark.sql.functions import lit
from pyspark.sql.types import StringType

# Add the "id" field with a default value if it is missing
if "id" not in grp_date_insights.columns:
    |   grp_date_insights = grp_date_insights.withColumn("id", lit("default_id"))

# Convert the "id" field to string if it is not already
if grp_date_insights.schema["id"].dataType != StringType():
    |   grp_date_insights = grp_date_insights.withColumn("id", grp_date_insights["id"].cast(StringType()))

# Save the DataFrame
grp_date_insights.write.format("cosmos.oltp").options(**config).mode("append").save()
```

▶ (14) Spark Jobs