Azure Databricks with Azure Data Factory
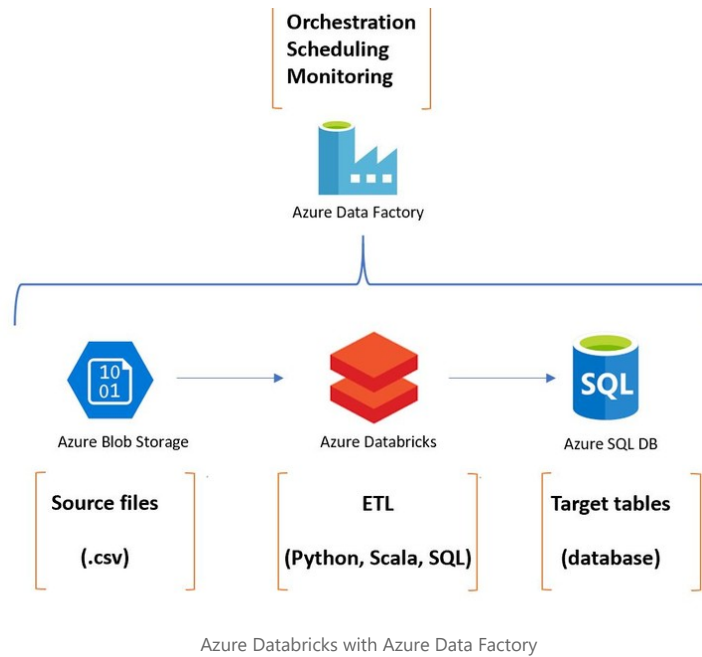
# ETL Modernization using Azure Databricks and Azure Data Factory

Published on February 26, 2019

**Rahul Athale**
Director, Azure Analytics & AI, Customer Success - Worldwide Commercial Business

**6 articles**    **+ Follow**

A simple example depicting commonly used ETL activities like: (1) Joining Two Datasets (2) Renaming Column (3) Creating a new column with a computed value.

## Source Data:



Two .csv files on Azure Blob - "Sales.csv" & "Product.csv"

## ETL Activities:

Use an Azure Databricks Python Notebook to perform these transformation tasks.

- *Mount an Azure Blob Storage container as DBFS*

Messaging

```
dbutils.fs.mount(
    source = "wasbs://"+getArgument("filepath")+"@"+storageName+".blob.c
    mount_point = "/mnt/adfdata",
    extra_configs = {"fs.azure.account.key."+storageName+".blob.core.win
                    accessKey})
```

- *Access files in your container as if they were local files*

```
%python

df_sales = sqlContext.read.format('csv').options(header='true', inferSch
df_products = sqlContext.read.format('csv').options(header='true', infer
```

- *Join two dataframes ; Rename column and sort ; Create ne*

```
from pyspark.sql.functions import col

# Join two dataframes on a column called "Product"
df_merged = df_sales.join(df_products, ["Product"])

# Rename column to new name, sort on one column
df_merged_renamed = (df_merged
                     .withColumnRenamed("Description", "Product_Description")
                     .sort("Transaction_date")
                    )

# create new column with Computed value and call it doublePrice
df_doubleprice = (df_merged_renamed
                  .withColumn("DoublePrice", col("Price") * 2)
                  .select("Product", "Transaction_date", "Price", "DoublePrice", "Product_Descriptic

display(df_doubleprice)
```

Sorted        Computed        renamed

▶ (1) Spark Jobs
▶ ▦ df_merged: pyspark.sql.dataframe.DataFrame = [Product: string, Transaction_date: string ... 11 more fields]
▶ ▦ df_merged_renamed: pyspark.sql.dataframe.DataFrame = [Product: string, Transaction_date: string ... 11 more fields]
▶ ▦ df_doubleprice: pyspark.sql.dataframe.DataFrame = [Product: string, Transaction_date: string ... 12 more fields]

| Product | Transaction_date | Price | DoublePrice | Product_Description | Payment_Type | Name | City | State | Country | Account_Created | Last_Login | Latitude | Longitude |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Product2 | 1/1/2009 10:06 | 3600 | 7200 | Product Two | Mastercard | Irene | Munich | Bayern | Germany | 1/1/2009 9:13 | 1/25/2009 11:37 | 48.15 | 11.583333 |
| Product1 | 1/1/2009 11:05 | 1200 | 2400 | Product One | Diners | Janis | Ballynora | Cork | Ireland | 12/10/2007 12:37 | 1/7/2009 1:52 | 51.8630556 | -8.58 |
| Product1 | 1/1/2009 12:19 | 1200 | 2400 | Product One | Visa | Marlene | Edgewood | WA | United States | 1/16/2006 14:45 | 1/17/2009 11:38 | 47.25028 | -122.2925 |
| Product2 | 1/1/2009 12:20 | 3600 | 7200 | Product Two | Visa | seemab | San Pawl tat-Targa | | Malta | 9/17/2005 3:32 | 1/22/2009 12:00 | 35.9202778 | 14.4425 |
| Product2 | 1/1/2009 12:25 | 3600 | 7200 | Product Two | Mastercard | Anne-line | Zug | Zug | Switzerland | 3/29/2005 23:14 | 1/31/2009 10:58 | 47.1666667 | 8.5166667 |
| Product1 | 1/1/2009 12:42 | 1200 | 2400 | Product One | Visa | ashton | Exeter | England | United Kingdom | 12/15/2008 1:16 | 2/9/2009 2:52 | 50.7 | -3.5333333 |

▦  ▦  📊 ▾  ⬇

Command took 4.41 seconds -- by raathale@microsoft.com at 2/11/2019, 10:31:46 PM on ETL_Cluster_Demo

## Target Table:
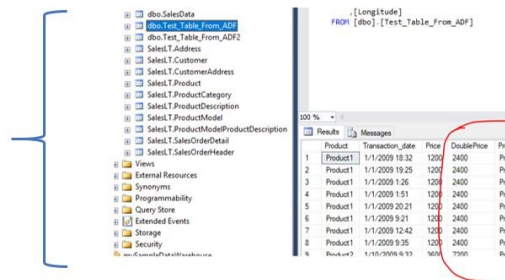
👍 Like   💬 Comment   ↪ Share                    👍 120 · 7 Comments

```
jdbcUrl = "jdbc:sqlserver://{0}:{1};database={2};encrypt=true;trustServe
```

- *Load/Write ETL output to Table in SQL Server*

```
db_target_table = getArgument("target_tablename")
df_doubleprice.write.jdbc(jdbcUrl, db_target_table, properties=connectio
```



**Target table on SQL DB**

A database table *(with name based on value of parameter "targe*
created or updated in the SQL Server database. This table is the result set of joining
"Sales.csv" & "Product.csv" + renaming of "Description" column + creation of new
computed column "DoublePrice".

## Orchestration using Azure Data Factory:

**Pipeline:**

- Lookup activity: to check for source data availability at directory location

- Copy activity: to copy processed files into different directory location

- Databricks Notebook activity: ETL tasks and load data into SQL Database

*Two parameters are passed from Azure Data Factory to Azure Databricks -
"filepath" & "target_tablename"*

Search Activities

▶ Batch Service
▶ Databricks
▶ Move & Transform
▶ Data Lake Analytics
▶ General
▶ HDInsight
▶ Iteration & Conditionals
▶ Machine Learning

Lookup
Availability Flag

Copy Data
copyFile

General | **Parameters** | Variables | Output

+ New | 🗑 Delete

| | NAME | TYPE | DEFAULT VALUE |
|---|---|---|---|
| | target_tablename | String ▼ | Test_Table_From_ |
| | filepath | String ▼ | testpath |

## Activity Runs:

## Activity Runs

Pipeline Run ID **9890d113-b2f0-4cdf-8995-dc88f2ad8fe5**

**All** | Succeeded | In Progress | Failed | Cancelled

| ACTIVITY NAME | ACTIVITY TYPE | ACTIONS | RUN START ⇕ | DURATION | STATUS | INT |
|---|---|---|---|---|---|---|
| Availability Fl... | Lookup | →] →] | 02/25/2019 1:33 PM | 00:00:01 | ✅ Succeeded | De |
| copyFile | Copy | →] →] 👓 | 02/25/2019 1:33 PM | 00:00:16 | ✅ Succeeded | De |
| ETL Activities | DatabricksNo... | →] →] | 02/25/2019 1:33 PM | 00:05:51 | ✅ Succeeded | DefaultIntegrationRuntime (East US) |

## Databricks Job:

**Microsoft Azure**

ADF_ADFwithDatabricks_pipeline1_ETL Activit

‹ All Jobs   View:   [ Code              ▼ ]      ☁ Export to

## ADF_ADFwithDatabricks_pipeline1

**Started:** 2019-02-25 13:33:55 CST
**Duration:** 5m 40s
**Status:** Succeeded
**Run ID:** 16
**Task:** Notebook at /Users/raathale@microsoft.com/etl_demo_python
  ▸  Parameters:
       {"filepath":"testpath","target_tablename":"Test_Table_From_ADF3","p
**Cluster:** Driver: Standard_DS3_v2, Workers: Standard_DS3_v2, 2 worke

## Output

### Define Widgets/ Parameters

Creating widgets for leveraging parameters to be passed

#Azure #AzureDatabricks #AzureDataFactory #Microsoft

Report this

---

## Published by

**Rahul Athale**
Director, Azure Analytics & AI, Customer Success - Worldwide Commercial
Business
Published • 1yr

6 articles   [ + Follow ]

**#Azure #AzureDatabricks #AzureDataFactory #Microsoft**

## Reactions

👤 👤 👤    👤 👤 👤 👤 👤 👤    (...)

## 7 Comments

Most Relevant ▼

[ Add a comment...                                              📷 ]

**Hayk 'Hike' C.** • 3rd+                                      1y  ...
Let's Talk Data | *Hiring*

How do you get all of the data into azure blob storage?

👍 · 1 Like  |  💬  1 Reply

Messaging

you can upload into Azure Blob using multiple options. That isn't the best

**Yago Luiz** • 3rd+
Senior Backend Developer at Wiz Soluções | Master's degree at University of Brasília
**Sidney Oliveira Raphael Brito**

👍 · 3 Likes 💬

**Load more comments**

**Rahul Athale**

Director, Azure Analytics & AI, Customer Success - Worldwide Commerce

**+ Follow**

n Rahul Athale

| | | | |
|---|---|---|---|
| tream Analytics & IoT | Azure Machine Learning – a simple example | Kafka + Spark Streaming + ADLS | Azure SQL DW - cost savings scalable DWUs |
| e on LinkedIn | Rahul Athale on LinkedIn | Rahul Athale on LinkedIn | Rahul Athale on LinkedIn |

icles

Messaging