

RAMP MLS Playbook

Created by Devash Ameta, Last modified by Chandan Bhendri on Apr 26, 2023

Introduction:

The purpose of RAMP MLS Playbook is to guide and help SAS users to migrate to other Applications i.e. to Pyspark, Python and Hive. The Table of Contents includes the prerequisites required to access certain data used to be available in RAMP, how to import the data in different applications, the platform limitations that users can face, how to make the connections with the different applications and most importantly, how to convert key SAS functionalities into other applications.

- > 1.) Limited CPU Utilization in python

> 2.) User Home Directory Space Limit in Cornerstone Gold

Cornerstone gold environment has limited space of 5GB in user home directory (</d1/home/user_ADS_Id>). If home directory quota limit exceeds, users will get an email like below:



Volume Details

Cluster Name	Volume Name	Mounted On	OS Disk Use	Quota Size	Usage
goldcluster	sysAdmin	/d1/home	9.0GB	5 GB	9.0GB

If you need more space, you can use your 'use case' path to import the data. Use Case Path Example: "`/expgrms/user_yarn_queue/devUser_ads_id`".

In above example, "user_yarn_queue" is the use case name. The Commands to find user specific use case name is documented in this MLS playbook Table of Contents header >> All users need is to replace the use case name in above example and keep the rest of the path same.

> 3.) How to download YB data into gold

Just like in RAMP SAS, users can still use ybunload function to download YB data in gold cluster. Using this ybunload function users can download the whole table from YB into Gold user use case. Below are the instructions:

- 1.) Open putty with the host name "`putty@edge.gso.aexp.com`" and Port "22".
- 2.) Login with your ads ID and password.
- 3.) Run the below ybunload - examples to download data from (YB) to gold cluster . Unloads the results of a query to a file csv format.

a.) To download the data from mina.yb to gold cluster, use below command:

```
ybunload -hd http://603a.gso.aexp.com:5402 -d ybprd02 -format csv -username <username> -W 4 mina.yb1_mer_dim -o /expgrms/your_own_yarn_queue/dev/your_ads_id/unload -prefix yb1_mer_dim -truncate-existing
```

b.) To download the data from CSIA(YB) to gold cluster, use below command:

```
/expgrms/ybload -A http://YB10a.gso.aexp.com:42100/unload -d csionetb3 -username <username> -W -format csv -o /expgrms/your_yarn_queue/dev/your_ads_id/unload -truncate-existing -select "select product_cd,lost_qty,ads_id from sparc_monthly_hist limit 10"
```

c.) Below query unloads a table from tmp schema:

```
/expgrms/ybload -A http://YB10a.gso.aexp.com:42100/unload -d csionetb3 -tmp_cbs_0502_all_email -username <username> -W -o /expgrms/your_yarn_queue/dev/your_ads_id/unload -truncate-existing
```

Note: keep everything else the same in above example except below:

username = Your ads ID in <username>,

Your_own_yarn_queue: Provide your specific use case name,

Your_ads_id = Keep your ads ID here.

d.) Establish Python or PySpark Connection using MLS (Access data from Cornerstone)

1.) Browse "`https://mlnbauto.gsd.aexp.com:443`".

2.) Select kernel as Python3.

3.) Click on Create and open notebook.

4.) Check your yarn quota by running the below commands in MLS or in putty respectively and change the same in below python/pyspark connection:

Step1: Run below command to check what CS UseCase(s) you have access:

```
Impred queue->showadcl | grep SUBMIT //this command to use in MLS?
```

```
mapred queue->showadcl | grep SUBMIT //just FYI on a separate note: this command to use in putty to check your Use Case(s)?
```

Step2: To run Python or PySpark query, you need to first use below code snippet to establish connection with CS by using your own Use Case,

*****Important note*****

Copy & paste below code block directly into the notebook. Make sure each and every single command line start from the very beginning of the cell as shown in below screenshot.

```
import os
import sys
import jdo
os.environ["SPARK_HOME"] = "/opt/maven/spark/spark"
os.environ["PYSPARK_PYTHON"] = "/usr/bin/python3.6"
spark_python = os.path.abspath(os.path.join(os.path.dirname(__file__), "python"))
sys.path.append(spark_python)
sys.path.append(os.path.join(spark_python, 'lib', 'py4j-0.10.7-src.zip'))
sys.path.append(os.path.join(spark_python, 'py4j'))
```

```
import os
import sys
import jdo
os.environ["SPARK_HOME"] = "/opt/maven/spark/spark"
os.environ["PYSPARK_PYTHON"] = "/usr/bin/python3.6"
spark_python = os.path.abspath(os.path.join(os.path.dirname(__file__), "python"))
sys.path.append(spark_python)
sys.path.append(os.path.join(spark_python, 'lib', 'py4j-0.10.7-src.zip'))
sys.path.append(os.path.join(spark_python, 'py4j'))
```

```
sc.setPythonExecutors(1)
sparkSession = SparkSession.builder.getOrCreate()
sqlContext = SparkSession.builder
```

```
.appName("Python Spark SQL Basic Example")
.config("spark.master", "yarn")
.config("spark.yarn.queue", "adls")  
----- use your yarn queue here (i.e., use case name to make the connection)
.config("spark.speculation", "true")
.config("spark.reducer.maxSizeInFlight", "1024")
.config("spark.jars", "/tmp/platform/disk/bulksparkjar/expPlatformLib/app/sparkhive/contrib.jar")
.config("spark.driver.maxMemory", "512m")
.config("spark.executor.instances", "20")
.config("spark.executor.memory", "2G")
.config("spark.executor.cores", "2")
.config("spark.executor.temporaryStorage.enabled", "true")
.config("spark.yarn.jars", "expPlatform/disk/bulksparkjar/expPlatformLib/app/sparkhive/contrib.jar")
.config("spark.yarn.queue", "adls")
.config("spark.submit.pythonFormat", "true")
.config("spark.submit.pythonArgs", "expPlatform/disk/bulksparkjar/expPlatformLib/app/sparkhive/contrib.jar")
.config("spark.submit.pythonPath", "20000")
.config("spark.submit.pythonVersion", "3.6")
.config("spark.yarnSupport", "true")
.config("spark.yarnContainerSupport", "true")
```

```
Raise RuntimeError(sqlContext.sparkContext, object)
dc = sqlContext.sparkContext
```

```
print(dc)
print(dc.context)
```

Note:
Users can add or change below config parameters in the above connection code as per their requirements to improve the performance of the whole code. In below config parameters we are using 4g as just an example for driver and executor memory. Users can change executor and driver memory and instances of executor as per their need.

```
config("spark.executor.memory", "4g"),
config("spark.executor.coreMemory", "10g"),
config("spark.driver.memory", "4g")
```

➤ 5.1) How to Create a Database?
In Order to save your data in your own database, you need to create database first like below by giving your own path in python/pyspark:

```
sqlContext.sql("CREATE DATABASE IF NOT EXISTS joins LOCATION '/mprfs/exp/gms/your_yarn_queue/devyouds/mstid0'")  
----- Use your own yarn queue (i.e., Use Case Name) and adds Id to create database at your own use case folder.
```

Below code is to create a Database in Hive:
CREATE DATABASE IF NOT EXISTS joins LOCATION '/mprfs/exp/gms/your_yarn_queue/devyouds/mstid0';
----- Use your own yarn queue (i.e., Use Case Name) and adds Id to create database at your own use case folder

How to use your own Hive database before running any code:
Use joins:

```
> 6.1) Cornerstone Package
from mtcgs import Cs3
obj = Cs3()
d = {}
try:
    "r": "/d/home/your_ads_id/hive",
    "o": "/d/home/your_ads_id/cs3_out",
    "md5_log_level": "ERROR"
except:
    out = obj.execute("d")
```

Users don't have to run Pyspark connection code to make the connection with the cornerstone if they use above cornerstone package.
Save your SqL query in filext and then run the above code in MLS.
After the above code ran successfully, Check your output file at the given "o" path.

> 7.) Establish Hive connection in Putty (Access data from Cornerstone)
1.) Open putty with the host name "ipofedge=edge.geo.mprfs.com" and Port "22".
2.) Login with your ads ID and password.
3.) Copy "/tmp/platform/disk/bulksparkhive" in putty by right click and press enter.

In Order to run Hive queries, we need to make database first and use it to run queries.
hive> CREATE DATABASE IF NOT EXISTS joins LOCATION '/m/home/your_ads_id/hive';
hive> USE joins;

> 8.) Establish Python or Pyspark Connection in Putty (Access data from Cornerstone)
1.) Open putty with the host name "ipofedge=edge.geo.mprfs.com" and Port "22".
2.) Login with your ads ID and password,
3.) Type Python3 in Putty and press enter,
4.) Check your yarn queue by running the below command in putty and change the same in below python connection code:

Step1: Run below command to check what CS UseCase(s) you have access:

mapped queue>ShowMeCS | grep SUBMIT

Step2: To run Python query, you need to first use below code snippet to establish connection with CS by using your own Use Case.

```
import os
import sys
os.environ["SPARK_HOME"] = "/opt/magnap/spark"
os.environ["PYSPARK_PYTHON"] = "/usr/bin/python3.6"
os.environ["PYSPARK_DRIVER_PYTHON"] = "/usr/bin/python3.6"
os.environ["PYSPARK_DRIVER_PYTHON_OPTS"] = "spark://your_spark_ip:7077"
py4j_glib_path = os.path.join(os.getenv("SPARK_HOME"), "None", "python")
py4j_glib = py4j_glib_path + "/py4j-0.10.7.jar"
sys.path.append(py4j_glib)
from pyspark import SparkSession
from pyspark.sql import SparkSession
sqlContext = SparkSession.builder \
    .appName("Python Spark SQL Basic Example") \
    .config("spark.yarn.queue", "admin") \
    .config("spark.submit.deployMode", "client") \
    .config("spark.yarn.queue", "admin") ----- use your yarn queue here (i.e., use case name to make the connection \
    .config("spark.yarn.access.hadoopFileSystems", "hdfs://your_spark_ip:8020") \
    .config("spark.sql.parquet.filterPushdown", "true") \
    .config("spark.hadoop.mapreduce.input.fileinputformat.input.dir.recursive", "true") \
    .config("spark.jars", "hadoop-mapreduce-client-core-2.12.1-hadoop-spark-3.2.1.jar") \
    .config("spark.executor.extraJavaOptions", "-XX:MaxDirectMemorySize=12m -XX:+UseConcMarkSweepGC -XX:+CMSInitiatingOccupancyFraction=30 -XX:+ScavengeBeforeFullGC -XX:+CMSParallelRemarkEnabled -XX:+UseCMSInitiatingOccupancyOnly -XX:+CMSParallelRemarkEnabled -XX:+UseCMSInitiatingOccupancyFraction=30 -XX:+ScavengeBeforeFullGC -XX:+CMSInitiatingOccupancyFraction=30 -XX:+ScavengeBeforeFullGC -XX:+CMSParallelRemarkEnabled") \
    .config("spark.yarn.executor.jars.hdfsPath", "hdfs://your_spark_ip:8020/user/your_spark_user/spark-earns-0.1.jar") \
    .config("spark.yarn.archive", "hdfs://your_spark_ip:8020/user/your_spark_user/spark-earns-0.1.jar") \
    .config("spark.yarn.jars", "hdfs://your_spark_ip:8020/user/your_spark_user/spark-earns-0.1.jar") \
    .config("spark.yarn.queue", "admin") \
    .config("spark.sql.execution.arrow.maxRecordsPerBatch", "20000") \
    .config("spark.sql.execution.arrow.enabled", "true") \
    .getOrCreate()
print(sqlContext)
print(sqlContext)
print(sqlContext.sparkContext)
print(sqlContext)
print(sqlContext.sparkContext.getOrCreate())
sc = sqlContext.sparkContext
print(sc)
```

Note:

Users can add or change below config parameters in the above connection code as per their requirements to improve the performance of the whole code. In the below config parameters we are using 4g as just an example for driver and executor memory. Users can change executor and driver memory and instances of executor as per their need.

```
.config("spark.executor.memory", "4g") \
.config("spark.executor.instances", "10") \
.config("spark.driver.memory", "4g")
```

sqlContext.set("CREATE DATABASE IF NOT EXISTS **jorts** location 'maprfs://**your_spark_ip**/**your_yarn_queue**/deviyouards/mStudio'"); *----- Use your own yarn queue (i.e., Use Case Name) and add id to create database at your own use case folder.*

> 9.) Execution time (Time function)

If we want to capture the execution time in SAS, MLS or Python.

Use the below code snippets at the start and end of the code to know the execution time.

Time Function:

	Time() in SAS	Time() in MLS	Time() in Python
At the start of the code:	%utlfunc(date()); %utlfunc(time());	Use %time at the start of the code.	import time start = time.time()
At the end of the code:	%utlfunc(date()); %utlfunc(time());	NA	end = time.time() print(end - start)

> 10.) How to import the data ?

How users can import the data is covered below at "Converting Key SAS Functions into Pyspark, Hive and Python"

> 11.) EXAMPLE! A User SAS Code End to End Conversion to Pyspark

RAMP User SAS Code End To End Conversion into Pyspark

> 12.) How to move CSV files among RAMP, CS Gold, YB and Your Local D Drive ?



How_to_move_files_among_C_drive_CS_Gold_...

> 13.) FAQs

Q1: Can users access MRA/YB data in MLS or Putty ? (MRA libraries: GMWDM, GMWDMBE, GMWDGMGP, GMWDMOPS, GMWDMPDS, GMWDMUDT, GMWDMPII, GMWDW, GMWDBE, GMWDWGMP, GMWDWPDS)?

Ans: NO, Right now users can access above MRA data only via RAMP. We are working with MRA and other teams to bring these data into CS, RAMP users will be promptly informed whenever we have updates in the progress. Stay tuned!

Q2: Can users do the joins between Yellowbrick and Cornerstone?

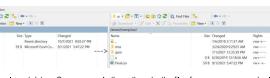
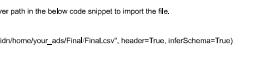
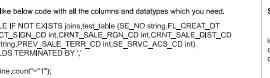
Ans: NO, RAMP team is working on this to get a solution. Will update once the solution is available.

Q3: Can users extract Yellowbrick data in MLS?

Ans: NO, RAMP team is working on this to get a solution. Will update once the solution is available.

Q4: Can users extract MCRM/DM/MCRM/WP data in MLS or in Putty?

Ans: NO, The data in these two libraries reside in MCRM database in YB. Right now users can only access these data via RAMP. However, we are working with MRA team to access the options to either direct access these data in YB via DBEAVER, and/or bring these data into CS, RAMP users will be promptly informed whenever we have updates in the progress. Stay tuned!

Note					
Query type	SAS	Pyspark	Hive	Python	
Import Data	<p>SAS Query</p> <p>Users can bring file into SAS servers in below ways:</p> <ul style="list-style-type: none">1. Via SAS EG : User can import the file from File option - Import File, give your required inputs and file will be automatically imported.2. Via WebSCP :Step1: Open WinSCPStep2: Drag your file from local to SAS server like below snapshot:  <p>Step3: Click On File--> Import data --> Servers --> go to the path under the files from your server and select your file.</p> <p>3. Via Code: Use the below code snippet to import the file:</p> <pre>PROC IMPORT DATAFILE ILE='\\dn\\home\\your_ads\\Final\\Final.csv'; OUT=Tracking_file DBMSSOCN REPLACE; RUN;</pre> <p>realtime 0.02 seconds cpu time 0.01 seconds</p>		<p>Pyspark query</p> <p>Step1: Open winscp.</p> <p>Step2: Drag your file from local to gold server like below snapshot:</p>  <p>Step3: Give your gold server path in the below code snippet to import the file.</p> <pre>df = sqlContext.read.csv('\\dn\\home\\your_ads\\Final\\Final.csv', header=True, inferSchema=True)</pre> <p>CPU times: user 11.2 ms, sys: 5.01 ms, total: 16.2 ms Wall time: 39 s</p>	<p>Hive Query</p> <p>Step1: Open winscp.</p> <p>Step2: Drag your file from local to gold server like below snapshot:</p>  <p>Step3: Create the external table like below with all the columns and datatypes which you need. <code>hive> CREATE EXTERNAL TABLE IF NOT EXISTS purchased_file (base_SE_NO string, LOAD_DATE string, TRIPUP_PACK_CD string, ACCTLSN_CD string, SALE_RQSL_CD string, CRNT_SALE_CD string, CRNT_TERL_CD string, PREV_SALE_TERL_CD string, SE_SVC_CD int) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' STORED AS TEXTFILE TBLPROPERTIES('skip.header.line.count'='1');</code></p> <p>Step4: Give your gold server path of the file in the below code snippet to import the data into the table.</p> <pre>hive> LOAD DATA LOCAL INPATH '\\dn\\home\\your_ads\\Final\\Final.csv' INTO TABLE joins.test_table;</pre> <p>Time taken: 15.8 seconds</p>	<p>Python Query</p> <p>Step1: Open winscp.</p> <p>Step2: Drag your file from local to gold server like below snapshot:</p>  <p>Step3: Give your gold server path in the below code snippet to import the file.</p> <pre>import pandas as pd df = pd.read_csv('\\dn\\home\\your_ads\\Final\\Final.csv') df.head()</pre> <p>0.0043 seconds</p>
Left Join with IN Operator	<p>Proc SQL; Create Table Trans_B2B_202102 as Select * from subm_se10.bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from RAMPCG.GMS_MERCHANT_CHAR left join RAMPCG.gms_transaction b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo IN ('202103','202102','202101','202012','202011','202010','202009','202008','202007','202006','202005','202004'); quit;</p> <p>real time 13:14.63 cpu time 1:56.65</p>		<pre>sqlContext.sql("" Create Table joins.Trans_InOp as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo IN ('202103','202102','202101','202012','202011','202010','202009','202008','202007','202006','202005','202004'); ")</pre> <p>CPU times: user 112 ms, sys: 60.4 ms, total: 172 ms Wall time: 9min 10s</p>	<pre>Create Table joins.Trans_InOp as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo IN ('202103','202102','202101','202012','202011','202010','202009','202008','202007','202006','202005','202004');</pre> <p>Time Taken: 538.145 seconds</p>	<pre>sqlContext.sql("") Create Table joins.Trans_InOp as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo IN ('202103','202102','202101','202012','202011','202010','202009','202008','202007','202006','202005','202004');</pre> <p>454.32 Seconds</p>
Left Join	<p>Proc SQL; Create Table Trans_B2B_202102 as Select * from subm_se10.bse_subm_ctry_cd, bse_subm_sro_sys_id,bse_trans_am, bse_trans_dt,bse_trans_vt_mo from RAMPCG.GMS_MERCHANT_CHAR left join rampcg.gms_transaction b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202102'; quit;</p> <p>real time 2:09.69 cpu time 3:88 seconds</p>		<pre>sqlContext.sql("") Create Table joins.Trans_B2B_202102 as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202102'; ")</pre> <p>CPU times: user 36.3 ms, sys: 17 ms, total: 47.4 ms Wall time: 2min 10s</p>	<pre>Create Table joins.Trans_B2B_202102 as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202102'; ")</pre> <p>152.869 seconds</p>	<pre>sqlContext.sql("") Create Table joins.Trans_B2B_202102 as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202102'; ")</pre> <p>168.59 seconds</p>
Left Join	<p>Proc SQL; Create Table Trans_B2B_202101 as Select * from subm_se10.bse_subm_ctry_cd, bse_subm_sro_sys_id,bse_trans_am, bse_trans_dt,bse_trans_vt_mo from RAMPCG.GMS_MERCHANT_CHAR left join rampcg.gms_transaction b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202101'; quit;</p> <p>real time 3:01.68 cpu time 5:85 seconds</p>		<pre>sqlContext.sql("") Create Table joins.Trans_B2B_202101 as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202101'; ")</pre> <p>CPU times: user 22.7 ms, sys: 7.54 ms, total: 30.2 ms Wall time: 1min 33s</p>	<pre>Create Table joins.Trans_B2B_202101 as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202101'; ")</pre> <p>122.827 seconds</p>	<pre>sqlContext.sql("") Create Table joins.Trans_B2B_202101 as Select b.subm_se10, bse_subm_ctry_cd, bse_subm_sro_sys_id, bse_trans_am, bse_trans_dt,bse_trans_vt_mo from cataloged3.gms_TRANSACTION a left join cataloged3.gms_TRANSACTION b on a.ae10\$b.subm_se10 where bse_subm_ctry_cd = '036' and bse_subm_sro_sys_id != 'F21' and bse_trans_vt_mo != '202101'; ")</pre> <p>122.87 seconds</p>
Union All	<p>Proc SQL; Create Table Trans_B2B_R12_TMP as Select * from Trans_B2B_202102 union all Select * from Trans_B2B_202101 quit;</p> <p>real time 21.66 seconds cpu time 20.47 seconds</p>		<pre>sqlContext.sql("") Create Table joins.Trans_B2B_R12_TMP as Select * from joins.Trans_B2B_202102 union all Select * from joins.Trans_B2B_202101 ")</pre> <p>CPU times: user 48.1 ms, sys: 21.2 ms, total: 69.3 ms Wall time: 1min 28s</p>	<pre>Create Table joins.Trans_B2B_R12_TMP as Select * from joins.Trans_B2B_202102 union all Select * from joins.Trans_B2B_202101 ")</pre> <p>169.47 seconds</p>	<pre>sqlContext.sql("") Create Table joins.Trans_B2B_R12_TMP as Select * from joins.Trans_B2B_202102 union all Select * from joins.Trans_B2B_202101 ")</pre> <p>66.18 seconds</p>
Extract a data from a table	<p>Proc SQL; Create Table Trans_B2B_R12 as Select * from Trans_B2B_R12_TMP union all Select * from Trans_B2B_R12;</p>		<pre>sqlContext.sql("") Create Table joins.Trans_B2B_R12 as Select * from Trans_B2B_R12_TMP union all Select * from Trans_B2B_R12;</pre>	<pre>Create Table joins.Trans_B2B_R12 as Select * from Trans_B2B_R12_TMP union all Select * from Trans_B2B_R12;</pre>	<pre>sqlContext.sql("") Create Table joins.Trans_B2B_R12 as Select * from Trans_B2B_R12;</pre>

	real time 14.14 seconds cpu time 14.13 seconds	real time 14.14 seconds cpu time 14.13 seconds	real time 14.14 seconds cpu time 14.13 seconds
ifElse Logic and Case statement	<pre>data trans_B2B_R12v1: where se_trans_am <= SUBM_AM_DCM from joins.Trans_B2B_R12_TMP """) CPU times: user 7.08 ms, sys: 5.13 ms, total: 10.2 ms Wall time: 23.2 s</pre>	<pre>Select *, se_trans_am <= SUBM_AM_DCM from joins.Trans_B2B_R12_TMP """) CPU times: user 7.08 ms, sys: 5.13 ms, total: 10.2 ms Wall time: 23.2 s</pre>	<pre>57.255 seconds</pre>
	<pre>if SUBM_AM_DCM < 0 then ROC <- 1 else ROC <- 1; format ROC :0.0; if SUBM_AM_DCM > 0 then ROC_Band = "Positive"; else if SUBM_AM_DCM < -29.99 then ROC_Band = "AUD-99.99 to 0"; else if SUBM_AM_DCM < -199.99 then ROC_Band = "AUD-199.99 to -100"; else if SUBM_AM_DCM < -199.99 then ROC_Band = "AUD-200.00+"; run;</pre>	<pre>if SUBM_AM_DCM < 0 then ROC <- 1 else ROC <- 1; format ROC :0.0; if SUBM_AM_DCM > 0 then ROC_Band = "Positive"; else if SUBM_AM_DCM < -29.99 then ROC_Band = "AUD-99.99 to 0"; else if SUBM_AM_DCM < -199.99 then ROC_Band = "AUD-199.99 to -100"; else if SUBM_AM_DCM < -199.99 then ROC_Band = "AUD-200.00+"; CPU times: user 38 ms, sys: 5.58 ms, total: 44.6 ms Wall time: 1min 16 s</pre>	<pre>create table Joins.Trans_B2B_R12v1 as select * case when SUBM_AM_DCM < 0 then '+' else '' end as ROC; case when SUBM_AM_DCM >= 0 then 'Positive' when SUBM_AM_DCM > -99.99 and SUBM_AM_DCM < 0 then "AUD-99.99 to 0" when SUBM_AM_DCM > -199.99 and SUBM_AM_DCM < -99.99 then "AUD-199.99 to -100" when SUBM_AM_DCM < -199.99 then AUD-200.00+ end as ROC_Band from joins.Trans_B2B_R12 where se_trans_dt >= 2021-01-01 and se_trans_dt <= 2021-02-28</pre>
Group By	<pre>proc sql Create table B2B_1 as Select SUM(SUBM_AM_DCM) as DBV_LCY, SUM(ROC) as ROCs, ROC_Band From joins.trans_B2B_R12v1 Group by ROC_Band; quit;</pre>	<pre>proc sql Create table B2B_1 as Select SUM(SUBM_AM_DCM) as DBV_LCY, SUM(ROC) as ROCs, ROC_Band From joins.trans_B2B_R12v1 Group by ROC_Band""") CPU times: user 20.4 ms, sys: 5.56 ms, total: 24.9 ms Wall time: 1min 16 s</pre>	<pre>61.419 seconds</pre>
Extract the data with no missing values.	<pre>bname joins /*<n>/home/cbha*/; proc sql Create table joins.B2B_Final as Select * From B2B_1 where (DBV_LCY is not missing and ROCs is not missing); quit;</pre>	<pre>selectContext.sql/*Create table joins.B2B_Final as Select * from joins.B2B_1 where DBV_LCY is not NULL and ROCs is not NULL; **/</pre>	<pre>Create table joins.B2B_1 as Select SUM(SUBM_AM_DCM) as DBV_LCY, SUM(ROC) as ROCs, ROC_Band From joins.trans_B2B_R12v1 Group by ROC_Band; 30.779 seconds</pre>
	<pre>real time 1.64 seconds cpu time 1.58 seconds</pre>	<pre>CPU times: user 2.84 ms, sys: 372 ps, total: 3.21 ms Wall time: 525 s</pre>	<pre>52.99 seconds</pre>
Extract the data in a new database	<pre>bname test /*<n>/home/cbha*/; proc sql Create table Test.Trans_B2B_202102 as Select * From Trans_B2B_202102; quit;</pre>	<pre>selectContext.sql/* Create table Test.Trans_B2B_202102 as select * from joins.Trans_B2B_202102""") CPU times: user 14.5 ms, sys: 257 us, total: 14.7 ms Wall time: 5.54 s</pre>	<pre>Create table Test.Trans_B2B_202102 as Select * from joins.Trans_B2B_202102; 81.095 seconds</pre>
Export the data	<pre>proc export data = share.B2B_Final out='<n>/home/your_pds_id/final.xlsx' dbms = 'XLSX'; real time 0.69 seconds cpu time 0.61 seconds</pre>	<pre>dtsqContext.sql/*"select * from Test.B2B_Final"*/ dtsqlToPandas() import pandas as pd df.to_excel('<n>/home/cbha/final.xlsx') CPU times: user 38.2 ms, sys: 64.1 ms, total: 102 ms Wall time: 4.68 s</pre>	<pre>insert overwriting local directory '<n>/home/cbha/join4' stored as text; select * from Test.B2B_Final; Time taken: 17,447 seconds</pre>
Create a database using macro	<pre>bname test /*<n>/home/cbha*/; %let a=%test; proc sql Create table a as select * from rampos.gms_transaction where se_subm_ctry_cd = '036' and se_trans_vr_mo IN ('202103','202102'); quit;</pre> <p>PROCEDURE SQL used: real time 0.58,11 cpu time 0.15 seconds</p>	<pre>by; db = 'test'; location = '%macro(cbha)' selectContext.sql/*CREATE DATABASE IF NOT EXISTS (a) LOCATION '(b)db'/*format(a\$db\$location)) except Exception as e: print(e) print("create table a as select * from rampos.gms_transaction where se_subm_ctry_cd = '036' and se_trans_vr_mo IN ('202103','202102') limit 10";format(a\$db)) sqContext.sql/*"select * from (a).new_table"*/,format(a\$db),show(); CPU times: user 4.7 ms, sys: 163 us, total: 4.87 ms Wall time: 9.51 s</pre>	<pre>set hivewar db1 = cbha_cham; set hivewar location = %macro(cbha); CREATE DATABASE IF NOT EXISTS (a) LOCATION '\$(b)hivewar\$(db)' Create table a(hivewar_db1) as select * from rampos.gms_transaction where se_subm_ctry_cd = '036' and se_trans_vr_mo IN ('202103','202102') limit 10; Time Taken: 232.46 seconds</pre>
Create a table with a left join using macro	<pre>bname test/*<n>/home/cbha*/; %let a=%test; %let b=%RAMPOS; proc sql CREATE TABLE A as trans_inpos_pyspark AS (Select b.subm_se10, b.se_subm_ctry_cd, b.se_subm_src_sys_id, b.se_trans_am, b.se_trans_dt,b.se_trans_vr_mo from b left join a on a.a10=b.a10 where a.a10 <= 1036 and b.se_subm_src_sys_id = 'F21' and b.se_trans_vr_mo IN ('202103') and b.se_trans_dt >= '2021-01-01')/*format(a\$a1,b\$b2)*/ print("Table Created successfully") except Exception as e: print(e) d = %sqlContext.sql("select * from (j.trans_inpos_pyspark)",format(db1)) d.show();</pre>	<pre>B1#1="test" db2="testbed3" selectContext.ad/*CREATE TABLE A(trans_inpos_pyspark AS (Select b.subm_se10, b.se_subm_ctry_cd, b.se_subm_src_sys_id, b.se_trans_am, b.se_trans_dt,b.se_trans_vr_mo from b left join a on a.a10=b.a10 where a.a10 <= 1036 and b.se_subm_src_sys_id = 'F21' and b.se_trans_vr_mo IN ('202103') and b.se_trans_dt >= '2021-01-01')/*format(a\$a1,b\$b2)*/ print("Table Created successfully") except Exception as e: print(e) d = %Context.sql("select * from (j.trans_inpos_pyspark)",format(db1)) d.show();</pre>	<pre>real time 14.93 seconds cpu time 13.60 seconds</pre>

	<pre> sys: 14.8 ms, total: 38.1 ms Wall time: 1min 5s </pre>	<pre> from pyspark.sql.functions import * from pyspark.sql.types import StringType,BooleanType,DateType df_transformed = df_transformed.withColumn("se_trans_am",cast(StringType,"se_trans_am").withCollation("Latin1_CI_AS_CS_WS")) df_transformed = df_transformed.withColumn("se_subm_usr_id",cast(StringType,"se_subm_usr_id").withCollation("Latin1_CI_AS_CS_WS")) df_transformed.show() </pre>	<pre> set hivevar:db1=test; create table \$table as select * from transformed AS (SELECT *, upper(trim(concat_ws(":",se_trans_am, se_subm_usr_id, se_trans_dt, se_trans_tm, se_trans_ip, se_trans_new))) as se_trans_ipspark_new FROM \$hivevar:db1).trans_ipos_pyspark; Time taken: 127.178 seconds </pre>	<pre> from pyspark.sql.functions import * from pyspark.sql.types import StringType,BooleanType,DateType df_transformed = df_transformed.withColumn("se_trans_am",cast(StringType,"se_trans_am").withCollation("Latin1_CI_AS_CS_WS")) df_transformed = df_transformed.withColumn("se_subm_usr_id",cast(StringType,"se_subm_usr_id").withCollation("Latin1_CI_AS_CS_WS")) df_transformed.show() CPU times: user 339.0 ms, sys: 407.0 ms, total: 37.8 ms Wall time: 1.08 s </pre>	<pre> Time taken: 127.178 seconds </pre>	
Create View using macro	<pre> Ulikewise test "f1dn/home/rba". proc sql create view &_vw_trans_ipos_pyspark_new as select se_trans_dt, se_subm_usr_id, se_trans_ip, se_trans_am, se_trans_tm, se_trans_ipspark_new from f1_transformed where se_trans_ip=202103 order by se_subm_usr_id; quit; </pre>	<pre> db1>test try: _d_transformed.createOrReplaceTempView("f1_db1") sqlContext.sql("""CREATE TABLE IF NOT EXISTS `trans_ipos_pyspark_new` SELECT * FROM f1_db1 WHERE se_trans_ip=202103 order by se_subm_usr_id"""); except Exception as e: print(e) try: sqlContext.sql("""Create or replace view &_vw_trans_ipos_pyspark_new AS select se_trans_dt, se_subm_usr_id, se_trans_ip, se_trans_am, se_trans_tm from f1_transformed where se_trans_ip=202103 order by se_subm_usr_id"""); except: print("View created successfully") _d_view_ipos_pyspark.show() _d_view = sqlContext.sql("Select * from &_vw_trans_ipos_pyspark_new").first(db1) _d_view.show() CPU times: user 5.52 ms, sys: 1.24 ms, total: 6.76 ms Wall time: 3.47 s </pre>	<pre> set hivevar:db1=test; Create or replace view \$view as select * from transformed where se_trans_ip=202103 order by se_subm_usr_id; </pre>	<pre> set hivevar:db1=test; CREATE TABLE IF NOT EXISTS `trans_ipos_pyspark_new` AS select se_trans_dt, se_subm_usr_id, se_trans_ip, se_trans_am, se_trans_tm from \$viewvar:db1.trans_ipos_pyspark; </pre>	<pre> Time taken: 0.477 seconds </pre>	<pre> db1>test try: _d_transformed.createOrReplaceTempView("f1_db1") sqlContext.sql("""CREATE TABLE IF NOT EXISTS `trans_ipos_pyspark_new` SELECT * FROM f1_db1 WHERE se_trans_ip=202103 order by se_subm_usr_id"""); except Exception as e: print(e) try: sqlContext.sql("""Create or replace view &_vw_trans_ipos_pyspark_new AS select se_trans_dt, se_subm_usr_id, se_trans_ip, se_trans_am, se_trans_tm from f1_transformed where se_trans_ip=202103 order by se_subm_usr_id"""); except: print("View created successfully") _d_view_ipos_pyspark.show() _d_view = sqlContext.sql("Select * from &_vw_trans_ipos_pyspark_new").first(db1) _d_view.show() CPU times: user 5.52 ms, sys: 1.24 ms, total: 6.76 ms Wall time: 3.47 s </pre>
count distinct	<pre> proc sql create table d1_distinct_count as select count(distinct(subm_se10)) as subm_se10_new from vw_ipos_ipos_pyspark_new quit; </pre>	<pre> real time 15.15 seconds cpu time 15.20 seconds </pre>	<pre> from pyspark.sql.functions import countDistinct d1_distinct_Count = d1_view.select(countDistinct("subm_se10"),).toDF("subm_se10_new") d1_distinct_Count.show() CPU times: user 16.4 ms, sys: 5.91 ms, total: 22.7 ms Wall time: 15.4 s </pre>	<pre> SELECT COUNT(DISTINCT subm_se10) as subm_se10_new FROM \$viewvar:db1.vw_trans_ipos_pyspark_new; </pre>	<pre> Time taken: 0.161 seconds </pre>	<pre> from pyspark.sql.functions import countDistinct d1_distinct_Count = d1_viewedict(countDistinct("subm_se10"),).toDF("subm_se10_new") d1_distinct_Count.show() CPU times: user 16.4 ms, sys: 5.91 ms, total: 22.7 ms Wall time: 15.4 s </pre>
Proc freq	<pre> proc freq data=vw_trans_ipos_pyspark_new; table se_subm_usr_ipo_ipo_id missing; run; </pre>	<pre> real time 27.09 seconds cpu time 46.52 seconds </pre>	<pre> from pyspark.sql.functions import count frequencies = d1_view.groupBy("se_subm_usr_ipo_ipo_id").agg(count("se_subm_usr_ipo_ipo_id").alias("Frequency")) frequencies.show() #F1 Does not sum to 100% 100*sum(Frequency)/sum(Frequency) over(1) Percent' #f1_pct "sum(Frequency)/over(Frequency des Cumulative_Frequency." "sum(Percent) over(order by Frequency des Cumulative_Percent" frequencies.show() CPU times: user 3.15 ms, sys: 1.02 ms, total: 5.08 ms Wall time: 2.09 s </pre>	<pre> with f1 AS (SELECT x.se_subm_usr_ipo_ipo_id,x.occurances, 100*x.occurances/SUM(x.occurances) OVER() AS pcent FROM (SELECT se_subm_usr_ipo_ipo_id, count(se_subm_usr_ipo_ipo_id) as occurances FROM \$viewvar:db1.vw_trans_ipos_pyspark_new group by se_subm_usr_ipo_ipo_id) AS f1 OVER (PARTITION BY se_subm_usr_ipo_ipo_id ORDER BY f1.se_subm_usr_ipo_ipo_id) sum(f1.occurances) AS Cummulative_Frequency, sum(f1.pcent) over(order by f1.occurances desc) Cumulative_Percent FROM f1 </pre>	<pre> Time taken: 190.537 seconds </pre>	<pre> from pyspark.sql.functions import count frequencies = d1_view.groupBy("se_subm_usr_ipo_ipo_id").agg(count("se_subm_usr_ipo_ipo_id").alias("Frequency")) frequencies.show() #F1 Does not sum to 100% 100*sum(Frequency)/sum(Frequency) over(1) Percent' #f1_pct "sum(Frequency)/over(Frequency des Cumulative_Frequency." "sum(Percent) over(order by Frequency des Cumulative_Percent" frequencies.show() CPU times: user 3.15 ms, sys: 1.02 ms, total: 5.08 ms Wall time: 2.09 s </pre>
Proc means with N and N miss	<pre> proc means nmiss mean min max sum data=vw_trans_ipos_pyspark_new; run; </pre>	<pre> real time 21.62 seconds cpu time 44.64 seconds </pre>	<pre> from pyspark.sql.functions import avg mean_aggs = d1_view.select(avg("se_subm_usr_ipo_ipo_id").alias("mean")) mean_aggs.show() from pyspark.sql.functions import col if # F1 Rows cut off the rows which contain null, None, 'NULL', or empty string values and returns a count from a given column null_values = d1_view.filter("view_subm_se10.contains('None' null 'NULL' '')".size() #(col("subm_se10") == '' null 'NULL') size() #d1_view.select("view_subm_se10").size() #d1_view.select("view_subm_se10").size() #(col("subm_se10") == null 'NULL' '') size() #d1_view.select("view_subm_se10").size() #(col("subm_se10") == null 'NULL' '') size() # Total count minus the count of all null values gives the count of not null values. # If all the values are null, then the count of null values # data = [(null_val,0),(None_val,0),(Null_val,0)] #Creating a dataframe of null and not null counts #d_new = sqlContext.createDataFrame(data,['Null_val','Non_Null_val']) #d_new.show() CPU times: user 17.1 ms, sys: 4.07 ms, total: 21.2 ms Wall time: 8.42 s </pre>	<pre> selected_avg = col("se_subm_usr_ipo_ipo_id") as mean from \$viewvar:db1.vw_trans_ipos_pyspark_new; with table as (select mean when subm_se10 is null or '' is null then 0 when subm_se10 is null then 1 when subm_se10 is 'NULL' then 1 when subm_se10 = '' then 1 when subm_se10 is 'None' then 1 else 0 end) null_count, count as Total from \$viewvar:db1.vw_trans_ipos_pyspark_new select table.null_count as null_values, (table.mean - total.null_count) as Not_Null_values from table; </pre>	<pre> Time taken: 86.354 seconds </pre>	<pre> from pyspark.sql.functions import avg mean_aggs = d1_view.select(avg("se_subm_usr_ipo_ipo_id").alias("mean")) mean_aggs.show() from pyspark.sql.functions import col if # F1 Rows cut off the rows which contain null, None, 'NULL', or empty string values null_values = d1_view.filter("view_subm_se10.contains('None' null 'NULL' '')".size() #(col("subm_se10") == '' null 'NULL') size() #d1_view.select("view_subm_se10").size() #d1_view.select("view_subm_se10").size() #(col("subm_se10") == null 'NULL' '') size() #d1_view.select("view_subm_se10").size() #(col("subm_se10") == null 'NULL' '') size() # Total count minus the count of all null values gives the count of not null values. # If all the values are null, then the count of null values # data = [(null_val,0),(None_val,0),(Null_val,0)] #Creating a dataframe of null and not null counts #d_new = sqlContext.createDataFrame(data,['Null_val','Non_Null_val']) #d_new.show() CPU times: user 17.1 ms, sys: 4.07 ms, total: 21.2 ms Wall time: 8.42 s </pre>
Notkeyup with dupout	<pre> Data trans_ipos_pyspark2; set vw_trans_ipos_pyspark_new; run; proc sort nodupkey data=trans_ipos_pyspark2 dupout=df_duplicate; by subm_se10; run; </pre>	<pre> real time 1.037 seconds cpu time 1.031 seconds </pre>	<pre> from pyspark.sql.window import Window from pyspark.sql.functions import row_number # Averting a window function which will partition by the given columns and and order by the given column df_new = d1_view.partitionBy("subm_se10").orderBy("subm_se10") # Selecting all the rows which are duplicate and will be dropped from the main dataframe into a new dataframe df_duplicate = df_new.withColumn("row_number",row_number().over(w2).filter("row_number > 1")).drop("row") df_new.show() CPU times: user 33.9 ms, sys: 4.07 ms, total: 37.9 ms Wall time: 1.03 s </pre>	<pre> -> (rowNumber() + 1) is a regex which means selecting every column from the table except the row number column df_new = d1_view because it is used only for the transformation and has no significance in the result. # Averting a window function which will partition by the given columns and and order by the given column select x.* from(select x.* from(select * from(select * from(select * from(select * over(partition by subm_se10 order by subm_se10) as m from \$viewvar:db1.vw_trans_ipos_pyspark_new where x.uniq_x > 1) as t) as t1) as t2) as t3 </pre>	<pre> Time taken: 127.178 seconds </pre>	<pre> from pyspark.sql.window import Window from pyspark.sql.functions import row_number # Averting a window function which will partition by the given columns and and order by the given column df_new = d1_view.partitionBy("subm_se10").orderBy("subm_se10") # Selecting all the rows which are duplicate and will be dropped from the main dataframe into a new dataframe df_duplicate = df_new.withColumn("row_number",row_number().over(w2).filter("row_number > 1")).drop("row") df_new.show() CPU times: user 33.9 ms, sys: 4.07 ms, total: 37.9 ms Wall time: 1.03 s </pre>

		real time 19.44 seconds cpu time 4.620 seconds	# creating all the rows which are unique into a new dataframes. df_Lines = spark.createDataFrame(df.select("row_number").over(2).filter(col("row") == 1).drop("row")) df_Lines.show() df_Lines.show()	select '_c_)?)*"; from (select * from df where row_number() over(partition by subn_se10 order by subn_se10) as m from \$hivevar\$db).vw_trans_inpos_pyspark_new e1x where x.m = 1	# creating all the rows which are unique into a new dataframes. df_Lines = spark.createDataFrame(df.select("row_number").over(2).filter(col("row") == 1).drop("row")) df_Lines.show() df_Lines.show()	No labels
Sum(), min(), max(), Count(*)	proc set create table df_operations as select sum(se_trans_am) as total_Sum, count(se_trans_am) as total_Count, min(se_trans_am) as min, max(se_trans_am) as max, min(maxDate(se_trans_dt)) as min_date_symmdd10,, max(maxDate(se_trans_dt)) as max_date_symmdd10,, max(PU(se_trans_dt,ymmd10)) as max_date_formal_symmdd10,, from trans_inpos_pyspark; quit;	CPU times: user 4.69 ms, sys: 2.93 ms, total: 7.61 ms, Wall time: 5.9 s	# creating all the rows which are unique into a new dataframes. df_operations = sqlContext.sql("with t0 as (select minse_trans_dt) as min_date, max(se_trans_dt) as max_date, sum(se_trans_am) as total_Count, min(se_trans_am) as min, max(se_trans_am) as max, min(maxDate(se_trans_dt)) as min_date_symmdd10,, max(maxDate(se_trans_dt)) as max_date_symmdd10,, max(PU(se_trans_dt,ymmd10)) as max_date_formal_symmdd10,, from df_operations) select t0,min,max,sum,t0.total_Count,from_unixtime(unix_timestamp(t0.min_date,'yyyy-MM-dd'),'ddMMyyyy') as min_date_from_t0,from_unixtime(unix_timestamp(t0.max_date,'yyyy-MM-dd'),'ddMMyyyy') as max_date_from_t0 df_operations.show()	Time taken: 122.11 seconds Time taken: 11622.2 seconds	224 seconds	
	real time 0.68 seconds cpu time 0.68 seconds	CPU times: user 2.19 ms, sys: 490 µs, total: 2.68 ms, Wall time: 526 ms				1.34 seconds
Drop view	proc sql drop view vw_trans_inpos_pyspark_new quit;	real time 0.00 seconds cpu time 0.00 seconds	try: sqlContext.sql("drop view if exists vw_trans_inpos_pyspark_new");\$form(db1) print("View dropped successfully") except Exception as e: print(e)	drop view \$hivevar\$db).vw_trans_inpos_prepare_new;	try: sqlContext.sql("drop view if exists vw_trans_inpos_pyspark_new");\$form(db1) print("View dropped successfully") except Exception as e: print(e)	0.001 seconds
Delete an observation	proc set delete from trans_inpos_pyspark where se_trans_dt=20210401 quit;	real time 2.0357 cpu time 22.23 seconds	CPU times: user 1.62 ms, sys: 650 µs, total: 2.46 ms, Wall time: 235 ms	try: d = df.filter(se_trans_dt == "2021-04-01") d.createOrReplaceTempView("NewTable") sqlContext.sql("INSERT OVERWRITE TABLE \$trans_inpos_pyspark select * from NewTable");\$form(db1) print("rows deleted successfully") except Exception as e: print(e)	Create table \$hivevar\$db).trans_inpos_pyspark_temp \$hivevar\$db).trans_inpos_pyspark; insert into \$hivevar\$db).trans_inpos_pyspark_temp select * FROM \$hivevar\$db).trans_inpos_pyspark where se_trans_dt < '2021-04-01'; INSERT OVERWRITE TABLE \$hivevar\$db).trans_inpos_pyspark select * from \$hivevar\$db).trans_inpos_pyspark_temp drop Table \$hivevar\$db).trans_inpos_pyspark_temp;	try: d = df.filter(se_trans_dt == "2021-04-01") d.createOrReplaceTempView("NewTable") sqlContext.sql("INSERT OVERWRITE TABLE \$trans_inpos_pyspark select * from NewTable");\$form(db1) print("rows deleted successfully") except Exception as e: print(e)
Append the data	proc set CREATE TABLE chand_trans_inpos_pyspark AS (select a.se_subm_se10,a.se_subm_ctry_cd,a.se_subm_socn_sy1,a.se_trans_am, a.se_trans_dt,a.se_trans_vno from \$trans_inpos_pyspark left join &GMS_MERCHANT_CHAR b on a.se_subm_se10=b.se_subm_se10 where a.se_subm_se10 > '2020-01-01' and a.se_trans_dt > '2020-01-01' and a.se_subm_socn_sy1 != 'P21' and a.se_trans_vno IN ('2020102'); quit;	real time 8.92 seconds cpu time 8.77 seconds	CPU times: user 3.21 ms, sys: 1.09 ms, total: 4.3 ms, Wall time: 78.2 ms	d2 = sqlContext.sql("Select a.se_subm_se10,a.se_subm_ctry_cd,a.se_subm_socn_sy1,a.se_trans_am, a.se_trans_dt,a.se_trans_vno from \$trans_inpos_pyspark left join &GMS_MERCHANT_CHAR b on a.se_subm_se10=a.se_subm_se10 where a.se_subm_se10 > '2020-01-01' and a.se_trans_dt > '2020-01-01' and a.se_subm_socn_sy1 != 'P21' and a.se_trans_vno IN ('2020102')) d2.union(d2) d2.show()	CPU times: user 28.9 ms, sys: 13.1 ms, total: 42 ms, Wall time: 1min 34s	CPU times: user 26.9 ms, sys: 13.1 ms, total: 42 ms, Wall time: 1min 34s
Drop Tables	proc sql drop table trans_inpos_pyspark quit;	real time 1.91 seconds cpu time 1.91 seconds	try: sqlContext.sql("drop table \$trans_inpos_pyspark");\$form(db1) print("table deleted") except Exception as e: print(e)	DROP TABLE \$hivevar\$db).trans_inpos_pyspark;	try: sqlContext.sql("drop table \$trans_inpos_pyspark");\$form(db1) print("table deleted") except Exception as e: print(e)	0.099 seconds

4 Comments

 Charles A Flynt

PYI - exporting to PDF doesn't provide readable output. Can someone fix that? Be good to have a version of this in PDF but current export option does not allow for horizontal scrolling, so anything that doesn't fit on one page is lost. This means the side by side conversion examples, the most valuable part of the document, do not translate to the PDF version.

 Julia J Bao

this book was directly built in confluence due to the 3 programming tools comparison, PDF format not working, you will need to go to this confluence to read it

 Charles A Flynt

Thanks for the reply,

I would note, it's actually difficult to read in this confluence format too, even on large monitors, due to having to scroll horizontally. Would be better in another format, which would also lend itself to better conversion to PDF as a side benefit. Hopefully someone will revise it someday to be more accessible. It seems a good resource, just unfortunately not well formatted.

 Julia J Bao

Thanks for your feedback! Since RAMP has been decommissioned, SAS is no longer a programming language. So users like you would need to copy and paste the areas needed into own document. I completely understand it is quit an effort for end users to do so .