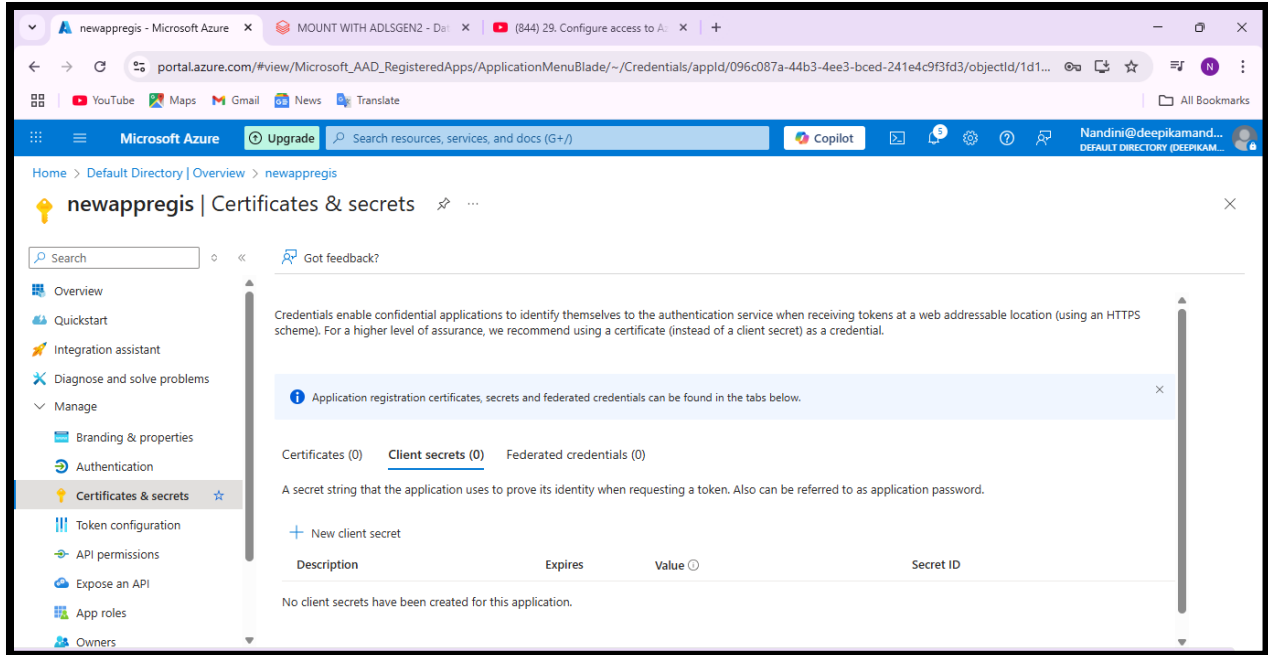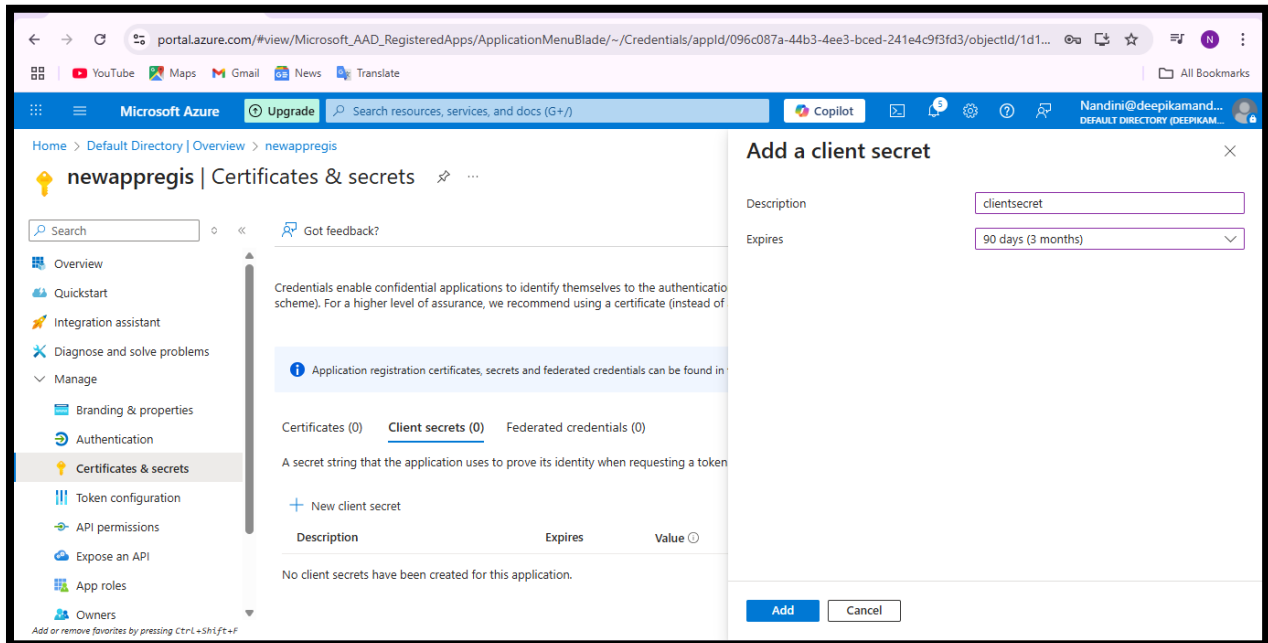# *PROJECT-06  (DATABRICKS)*

## **NANDINI RATHORE**

1.Python fundamentals: Data types, looping concepts (for, if), functions concept, exception handling (3 examples with project).

2.  Get data from ADLSGen2 and load in sql data base (service principle) (key vault) (Storage data contributor role, create scope)

3.  Use dynamic method to get data from sql database to adls gen2 (through url copy from another notebook)
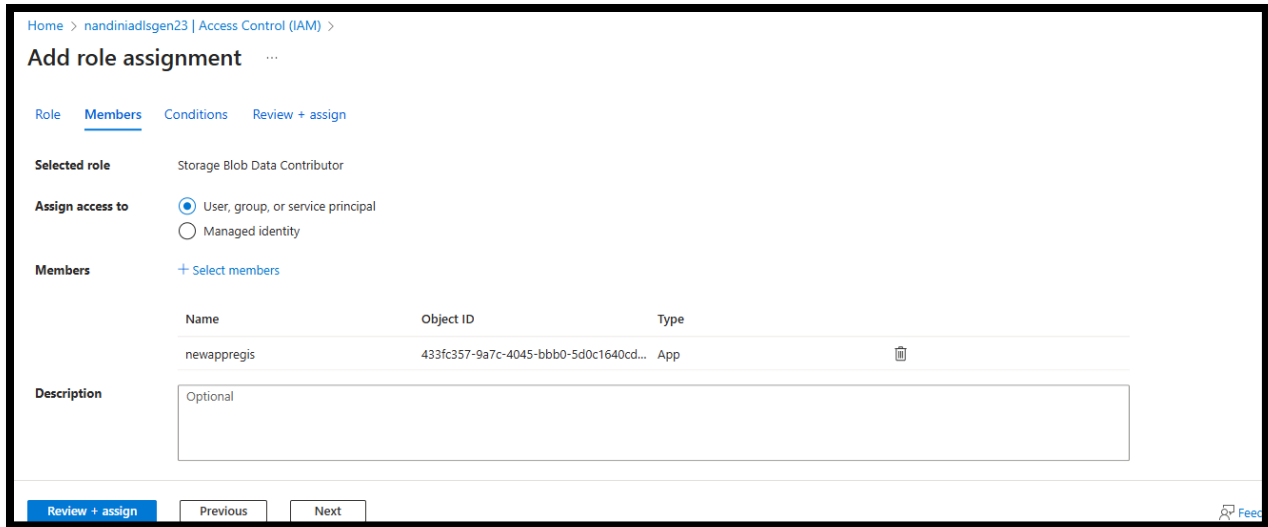
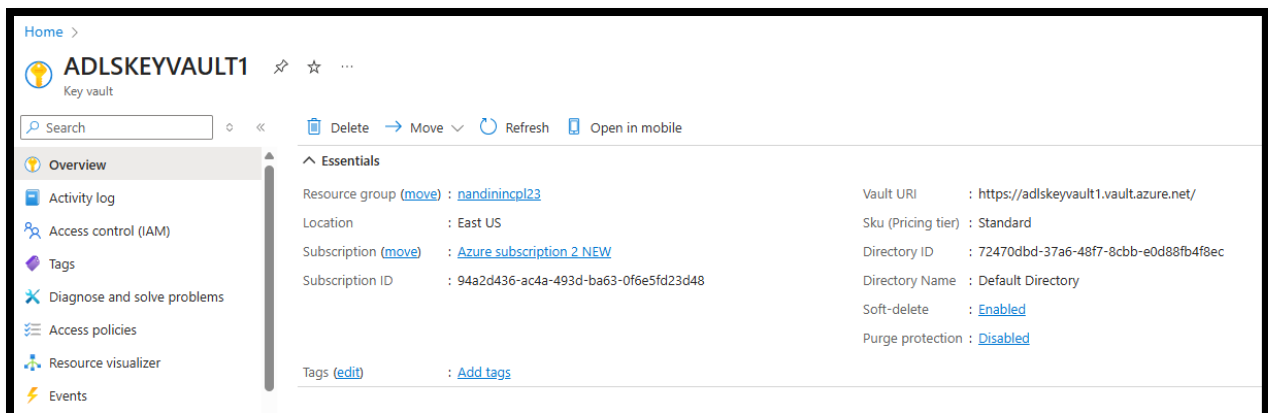1. OPEN AZURE→MICROSOFTENTRAID→APP REGISTRATION.

CREATE A NEW APP.

THEN GO TO CERTIFICATES AND SECRETS IN LEFT BAR AND CREATE CLIENT SECRETS AND COPY THAT SAVE IT SOMEWHERE IN DEVICE.



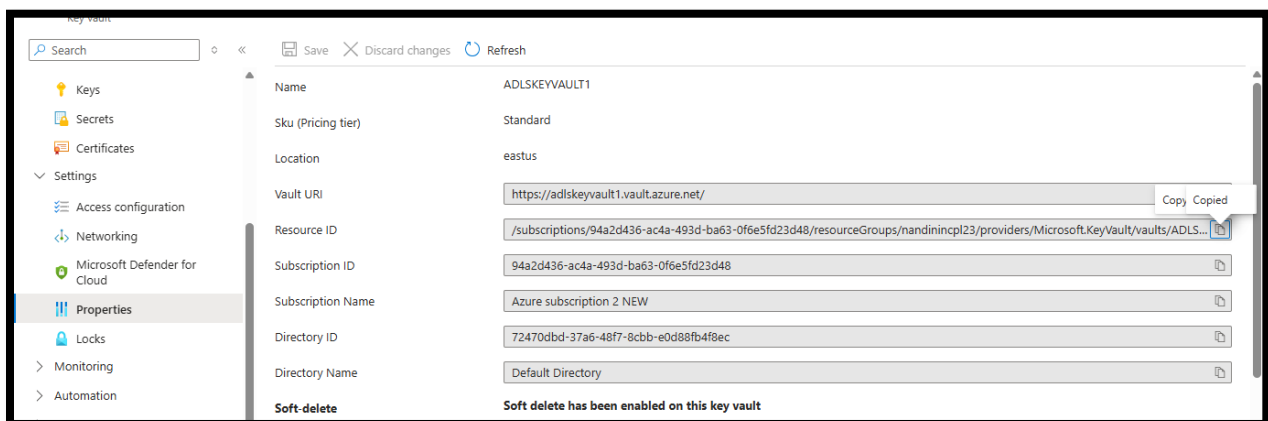2. CREATE CLIENT SECRET SELECT EXPIRYDATE.

3.  THEN ADLSGEN ACCOUNT→GIVE ROLE ASSIGNMENT TO REGISTERAPP→"STORAGE BLOB DATA CONTRIBUTOR."
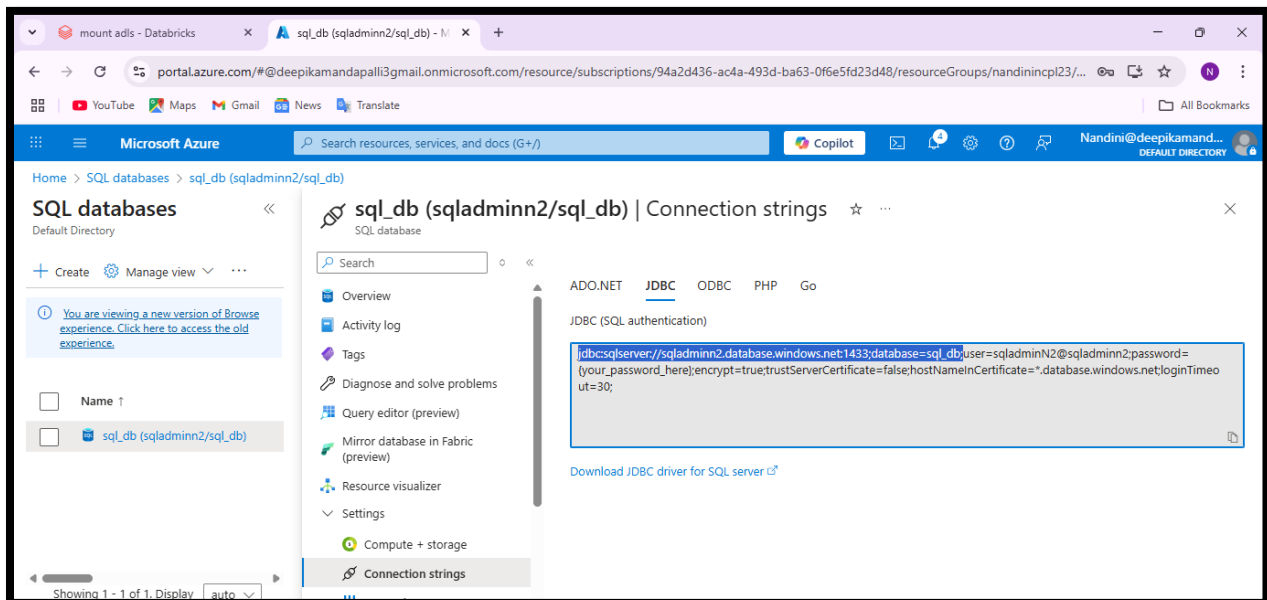


4.  FROM KEYVAULT,COPY DIRECTORY ID OR TENANT ID.

COPY VAULT URI-IT IS THE SPECIPIC ID OF KEY VAULT.

COPY RESOURCE ID .

CREATE SCOPE, WRITE NAME,PASTE VAULT URI IN DNS NAME AND PASTE RESOURCE ID THAT WE HAVE COPY FROM KEY VAULT.

5.mount with adls storage account.

Using scope ,to get clientid and client secret value.

Then mount the adls with databricks by using scope, clientid, clientsecret , storage container,storage account and mount point.



Abfss:→ Azure Blob File System Secure.

abfss:// is used to securely connect to Azure Data Lake Storage Gen2 from tools like Databricks and Spark.

List of scopes in databricks.



Secrets in the scope. Username and password used in sqldb.

## 1.DATA:- INCREMENT.MEDICAL_REPORTS1

## READ DATA FROM ADLSGEN2 TO SQLDB.

```
▶  ✓  ✓  4 days ago (2s)                                     4                          Python  ✦  []  ⋮

df=(spark.read
  .format("jdbc")
  .option("url", "jdbc:sqlserver://sqladminn2.database.windows.net:1433;database=sql_db;")
  .option("dbtable", "increment.MEDICAL_REPORTS1")
  .option("user", "sqladminN2")
  .option("password", "Ybcfp5xe7k")
  .load()
)
display(df)
```

▶ (1) Spark Jobs

▶ 🖿 df: pyspark.sql.dataframe.DataFrame = [R_ID: integer, PATIENT_ID: integer ... 2 more fields]
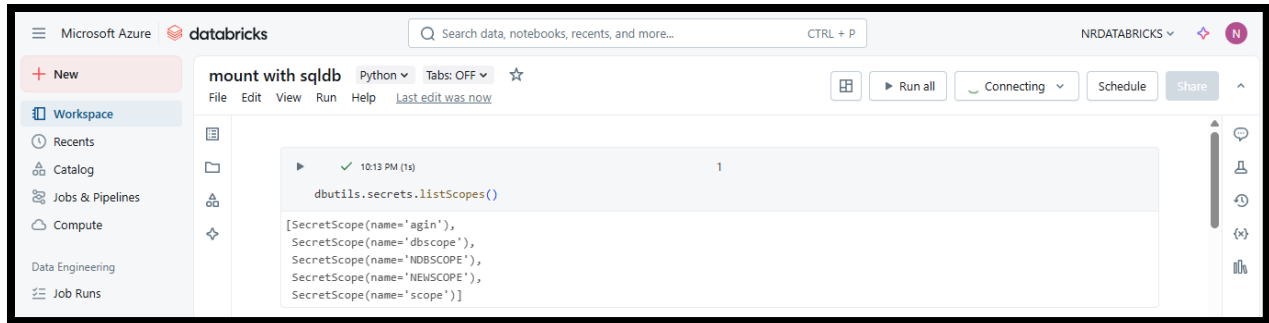
Table ∨     +                                                              Q ▽ 𝄚 ▢

| | R_ID | PATIENT_ID | R_TYPE | GENERATE_DATE |
|---|---|---|---|---|
| 1 | 10 | 1 | SUGAR TEST | 2025-06-28T12:00:00.000+00:... |
| 2 | 111 | 2 | THYROID | 2025-06-29T12:00:00.000+00:... |
| 3 | 222 | 3 | VITB12 | 2025-07-10T12:00:00.000+00:... |
| 4 | 333 | 4 | BLOOD TEST | 2025-06-29T12:00:00.000+00:... |

JDBC (Java Database Connectivity) is a connecting string to connect with database. Used in spark, databricks, java, scala.

## IN OTHER NOTEBOOK **DEFINE FUNCTIONS**.

```
▶  ✓  ✓  07:14 PM (<1s)                                     1                          Python  ✦  []  ⋮   🗑

class Exceptionraise(Exception):
    pass
```

```
▶  ✓  ✓  08:04 PM (2s)                                     2                          Python  ✦  []  ⋮   🗑

df=(spark.read
  .format("jdbc")
  .option("url", "jdbc:sqlserver://sqladminn2.database.windows.net:1433;database=sql_db;")
  .option("dbtable", "increment.MEDICAL_REPORTS1")
  .option("user", "sqladminN2")
  .option("password", "Ybcfp5xe7k")
  .load()
)
# TO CHECK TEST DATE OF FIRST FIVE PATIENTS.
def PATIENT_info(PATIENT_ID,R_TYPE,GENERATE_DATE):
        if int(PATIENT_ID) > 5:
            raise Exception("PATIENT_ID must be less than or equal to 5")
        else:
            print(f" PATIENT{PATIENT_ID} HAVE TEST {R_TYPE} ON {GENERATE_DATE} .")
```

CREATE NEW NOTEBOOK AND USE

**%run** is Databricks magic command to include another notebook.



This command tells Databricks to **run another notebook** located at the path.

- Any **functions, variables, or imports** defined in that notebook (NOTEBOOK)will become available in your **current notebook**.
- It's like copy-pasting all the code from NOTEBOOK into your current notebook — but cleaner and modular.

Call function and exception handling:

```python
PATIENT_ID=int(input("enter PATIENT_ID"))
row_found = df.filter(df['PATIENT_ID'] == PATIENT_ID).collect()
if row_found:
    row = row_found[0]
    try:
        PATIENT_info(row['PATIENT_ID'],row['R_TYPE'], row['GENERATE_DATE'])
    except Exception as e:
        print(e)
    finally:
        print("program completed")
else:
    print(f"no PATIENT found for THIS ID  {PATIENT_ID}")
```

▸ (1) Spark Jobs

enter PATIENT_ID

| 2 |
|---|

```
PATIENT2 HAVE TEST THYROID ON 2025-06-29 12:00:00 .
program completed
```

I define function in such a way that it ask to user for input and it define that which patient have test on which day.

If error occur then also it gives a mssg:-

```
PATIENT_ID=int(input("enter PATIENT_ID"))
row_found = df.filter(df['PATIENT_ID'] == PATIENT_ID).collect()
if row_found:
    row = row_found[0]
    try:
        PATIENT_info(row['PATIENT_ID'],row['R_TYPE'], row['GENERATE_DATE'])
    except Exception as e:
        print(e)
    finally:
        print("program completed")
else:
    print(f"no PATIENT found for THIS ID  {PATIENT_ID}")
```

▶ (1) Spark Jobs

enter PATIENT_ID

6

no PATIENT found for THIS ID  6

It gives an exception raise message " That no patient with id 6."

By using for loop:-

```
for row in df.collect():
    PATIENT_info(row.PATIENT_ID, row.R_TYPE, row.GENERATE_DATE)
```

▶ (1) Spark Jobs

```
PATIENT1 HAVE TEST SUGAR TEST ON 2025-06-28 12:00:00 .
PATIENT2 HAVE TEST THYROID ON 2025-06-29 12:00:00 .
PATIENT3 HAVE TEST VITB12 ON 2025-07-10 12:00:00 .
PATIENT4 HAVE TEST BLOOD TEST ON 2025-06-29 12:00:00 .
```

This command means it will "Go through each row in the Spark DataFrame df, and for each row, call the function patient_info()".

Write table from DATABRICKS TO ADLS

```
df.write.format("csv").save("/mnt/ncpl/csvfiless/PATIENTS.CSV")
```

▶ (1) Spark Jobs

Table stored in csv file in adls storage account.

**Second data:-STUDENT.FEE**

## 2nd example

```
✓ Just now (8s)                                          3

df=(spark.read
  .format("jdbc")
  .option("url", "jdbc:sqlserver://sqladminn2.database.windows.net:1433;database=sql_db;")
  .option("dbtable", "student.fee")
  .option("user", "sqladminN2")
  .option("password", "Ybcfp5xe7k")
  .load()
)
display(df)
```

Read table from sqldb in databricks.

| | $1^2_3$ s_rollno | $^A_C$ s_name | $^A_C$ course | $1^2_3$ FEE_PAID | $1^2_3$ DUE_AMOUNT | ⏱ LAST_PAID_DATE |
|---|---|---|---|---|---|---|
| 1 | 1 | Aditi | Python | 300 | 200 | 2024-11-01T00:00:00.000+00:... |
| 2 | 2 | Raj | Java | 500 | 0 | 2024-10-28T00:00:00.000+00:... |
| 3 | 3 | Meera | SQL | 200 | 300 | 2024-10-20T00:00:00.000+00:... |
| 4 | 4 | Arjun | Azure | 600 | 0 | 2024-11-05T00:00:00.000+00:... |
| 5 | 5 | Simran | Python | 100 | 400 | 2024-11-10T00:00:00.000+00:... |

5 rows | 7.66s runtime                                                    Refreshed now

```
✓ 12:05 AM (1s)                                          3                      Py

df=(spark.read
  .format("jdbc")
  .option("url", "jdbc:sqlserver://sqladminn2.database.windows.net:1433;database=sql_db;")
  .option("dbtable", "student.fee")
  .option("user", "sqladminN2")
  .option("password", "Ybcfp5xe7k")
  .load()
)
# TO CHECK WHICH STUDENT HAS DUE AMOUNT.
def check_due_fee(s_rollno, s_name, DUE_AMOUNT):
      if int(s_rollno) < 5:
          print(f"{s_name} (ROLL NO {s_rollno}) pending amount is ₹{DUE_AMOUNT}.")
      else:
          print(f"{s_name} (ROLL NO {s_rollno}) has cleared all dues.")

▸ 🗔 df: pyspark.sql.dataframe.DataFrame = [s_rollno: integer, s_name: string ... 4 more fields]
```

Define a function , to check which student fee is pending.



```
File  Edit  View  Run  Help    Python ⌄   Tabs: ON ⌄   ☆   Last edit was 1 minute ago                      ⊞   ▶ Run all   ● Nandini ⌄

  ▶        ✓ 1 minute ago (4s)                                          5

    s_rollno=int(input("enter s_rollno"))
    row_found = df.filter(df['s_rollno'] == s_rollno).collect()
    if row_found:
        row = row_found[0]
    try:
        if s_rollno <= 5:
            check_due_fee(row['s_rollno'],row['s_name'], row['DUE_AMOUNT'])
        else:
            print("you can see first 5 roll no.")
    except Exception as msg:
        print("msg")

    finally:
        print("program completed")

  ▶ (1) Spark Jobs

enter s_rollno

  7

you can see first 5 roll no.
program completed
```

Now call that function in other notebook having % run "path of previous notebook" and exception handling by using try ,except and final commands.

This screenshot is showing an **exception handling.**

```
    s_rollno=int(input("enter s_rollno"))
    row_found = df.filter(df['s_rollno'] == s_rollno).collect()
    if row_found:
        row = row_found[0]
    try:
        if s_rollno <= 5:
            check_due_fee(row['s_rollno'],row['s_name'], row['DUE_AMOUNT'])
        else:
            print("you can see first 5 roll no.")
    except Exception as msg:
        print("msg")

    finally:
        print("program completed")


▶ (1) Spark Jobs
```

enter s_rollno

3

Meera (ROLL NO 3) pending amount is ₹300.
program completed

It will ask for input and show result about that id.if that id is not present that raise an exception.

```
▶ ∨ ✓ Just now (3s)                              7                        Python  ✧ ⌗ ⋮
    df.write.mode("overwrite").format("csv").save("/mnt/ncpl/csvfiless/studentfee.csv")
▶ (1) Spark Jobs
```

Write table to adls as studentfee.csv file.

You can choose to store in csv ,parquet files as per requirements.

**3 Data:-Age_category_lookup:-**



```python
df=(spark.read
  .format("jdbc")
  .option("url", "jdbc:sqlserver://sqladminn2.database.windows.net:1433;database=sql_db;")
  .option("dbtable", "age_category_lookup")
  .option("user", "sqladminN2")
  .option("password", "Ybcfp5xe7k")
  .load()
)
display(df)
```

Read table age_category_lookup in table.

| | id | min_age | max_age | category_name | description |
|---|----|---------|---------|---------------|-------------|
| 1 | 1 | 0 | 12 | kid | child below teenage |
| 2 | 2 | 13 | 19 | Teen | Teenage phase |
| 3 | 3 | 20 | 35 | Young | Early adult life |
| 4 | 4 | 36 | 50 | Mid-Age | Mature adult life |
| 5 | 5 | 51 | 65 | Senior | Approaching retireme... |
| 6 | 6 | 66 | 120 | Elderly | Retired/Senior |

✓ 1 minute ago (1s)          7

```python
def age_category(min_age,max_age,category_name,description):
    if int(min_age)<0:
        raise Exception("age should not be less than 0.")
    elif int(min_age)>=0 & int(max_age)<12:
        print("kid")
    elif int(min_age)>=13 & int(max_age)<19:
        print("teen")
    elif int(min_age)>=20 & int(max_age)<34:
        print("young")
    elif int(min_age)>=35 & int(max_age)<50:
        print("middleage")
    elif int(min_age)>=51 & int(max_age)<65:
        print("senior")
    else:
        print("old")
```

▶ 🖿 df: pyspark.sql.dataframe.DataFrame = [id: integer, min_age: integer ... 3 more fields]

Define function and using if-else.

I use if-elif-else here because  it has multiple conditions.

```
▶  ∨  ✓ 1 minute ago (7s)                              9                          Python  ✧  [ ]  ⋮

    age=int(input("enter age"))
    row_found = df.filter((df['min_age']<=age) & (df['max_age']>=age)).collect()
    if row_found:
    |    row = row_found[0]
    try:
    |    age_category(row['min_age'],row['max_age'],row['category_name'],row['description'])
    except Exception as e:
    |    print(e)
    finally:
    |    print("program completed")


  ▸ (1) Spark Jobs

  enter age

  ┌──────────────────────────────────────────────────────────────────────────────────┐
  │ 15                                                                                 │
  └──────────────────────────────────────────────────────────────────────────────────┘

  kid
  program completed
```

Then calling function in another notebook and doing exception handling.

Row_found=df.filter((df['min_age']<=age) &(df['max_age']>=age)).collect()

**Each row** in the Spark DataFrame df, and for each row, **call the function** category_age().

.collect() :- Converts the whole Spark DataFrame into a list of Row objects .

**For row:-** Loop over each row[min_age & max_age] .

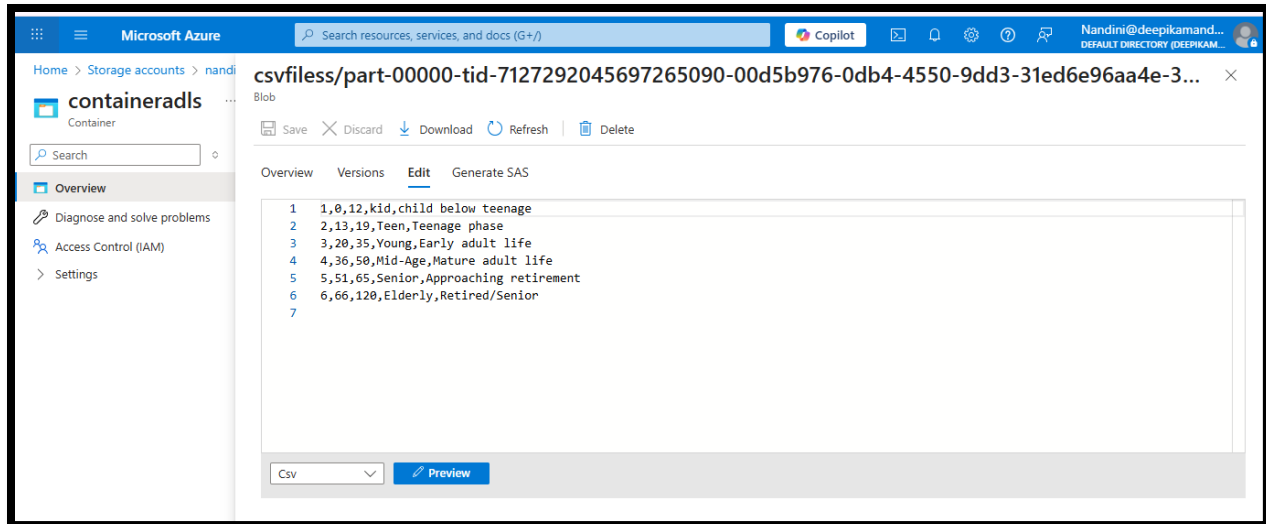Age_category  is a function  that we are calling it **once per row.**


## WRITE SQL TABLE IN ADLS IN CSV FILES.

```
▶  ∨  ✓ 2 minutes ago (5s)                              10                         Python  ✧  [ ]  ⋮

    df.write.format("csv").save("/mnt/ncpl/csvfiless")
  ▸ (1) Spark Jobs
```

Csv file in adls storage account.

## CONCLUSION:-

1. To ensure scalability, reusability, and maintainability, this project followed a modular approach by dividing the workflow into separate Databricks notebooks.
2. The first notebook was responsible for securely connecting to and reading data from Azure SQL Database using JDBC.
3. A second notebook was dedicated to defining reusable transformation functions, allowing for cleaner code and easier debugging.
4. The final notebook brought everything together by calling the other two notebooks. This method made the project more organized, easier to understand, and better for future updates. If anything needs to change later, we can do it in one place without changing the whole project.