# ASSIGNMENT-10

# (15ᵗʰ JUNE 2022)

NAME – MOHANA LIKHITHA THOTAKURA     ROLLNO – DXC-262AB-1219

BATCH – DXC-262-ANALYTICS-B12-AZURE     COMPANY – DXC TECHNOLOGY

EMPLOYEE DOMAIN – AZURE ANALYTICS     TRAINER NAME – MR. AJAY KUMAR

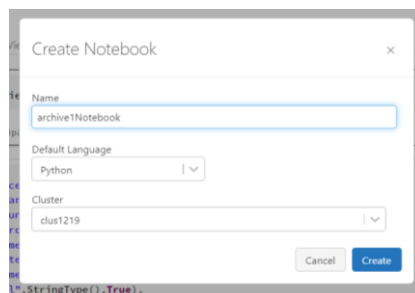DATE OF SUBMISSION – 15ₜₕ JUNE 2022     NO. OF QUESTIONS: 6

---

**Case 1**. Using archive1.zip file - please ingest data into Databricks DBFS path & query the data, redesign columns accordingly using Dataframe commands - display with notebooks accordingly.
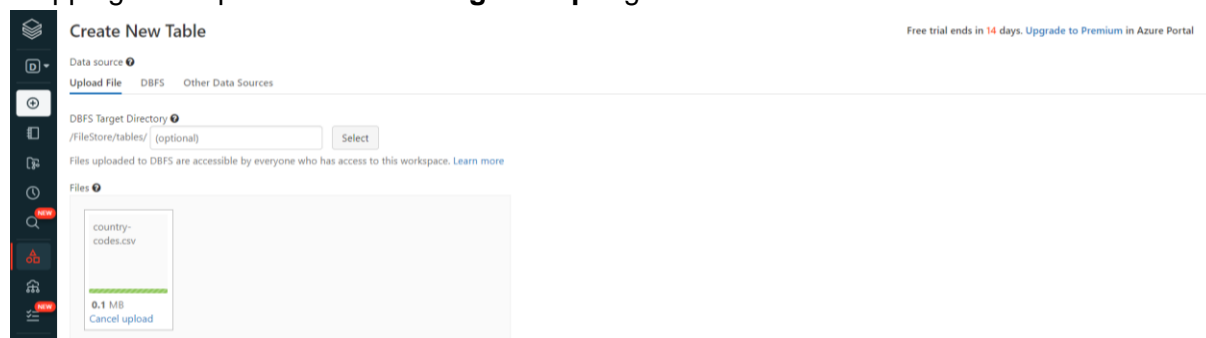**File Being used: country_codes.csv**
**Step 1:** First, login to your Azure Portal and create a Databricks workspace.
**Step 2:** Open the Databricks workspace and create cluster for your future use.
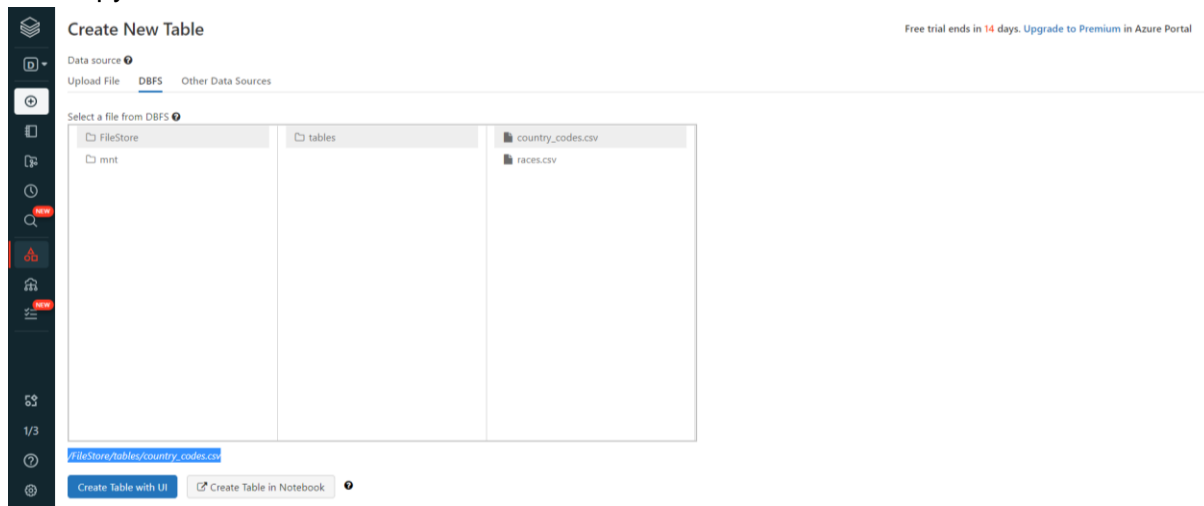**Step 3:** Now, create a notebook by clicking on the create Notebook option from the side panel.



**Step 4:** After creating the notebook, ingest the data into the Databricks by dragging and dropping the required file in the **drag & drop** region.
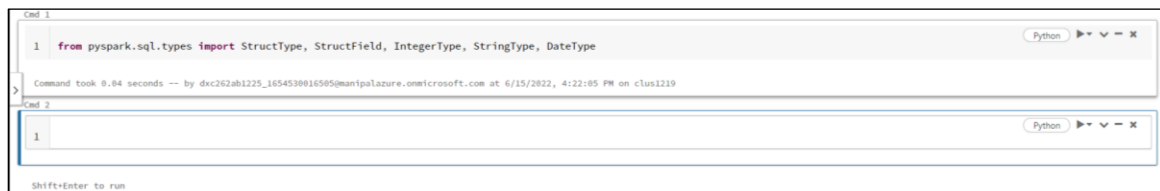
Later, click on DBFS and select the file that you have dropped. This will give you the file path and copy that.



**Step-5:** Import the required fields and features from pyspark.
*from pyspark.sql.types import StructType, StructField, IntegerType, StringType, DateType*



country_codes_schema = StructType(fields=[StructField("FIFA", StringType(),False),
StructField("Dial", StringType(),True),
StructField("ISO3166-1-Alpha-3",StringType(),True),
StructField("MARC", StringType(),True),
StructField("is_independent", StringType(),True),
StructField("ISO3166-1-numeric",IntegerType(),True),
StructField("GAUL", IntegerType(),True),
StructField("FIPS", StringType(),True),
StructField("WMO", StringType(),True),
StructField("ISO3166-1-Alpha-2",StringType(),True),
StructField("ITU", StringType(),True),
StructField("IOC", StringType(),True),
StructField("DS", StringType(),True),
StructField("UNTERM Spanish Formal", StringType(),True),
StructField("Global Code",StringType(),True),
StructField("Intermediate Region Code",IntegerType(),True),
StructField("official_name_fr",StringType(),True),
StructField("UNTERM French Short",StringType(),True),
StructField("ISO4217-currency_name",StringType(),True),
StructField("Developed / DevelopingCountries", StringType(),
True),

```
                                StructField("UNTERM Russian Formal",StringType(),True),
                                StructField("UNTERM English Short",StringType(),True),
                                StructField("ISO4217-
currency_alphabetic_code",StringType(),True),
                                StructField("Small Island Developing States
(SIDS)",StringType(),True),
                                StructField("UNTERM Spanish Short",StringType(),True),
                                StructField("ISO4217-
currency_numeric_code",IntegerType(),True),
                                StructField("UNTERM Chinese Formal",StringType(),True),
                                StructField("UNTERM French Formal",StringType(),True),
                                StructField("UNTERM Russian Short",StringType(),True),
                                StructField("M49",IntegerType(),True),
                                StructField("Sub-region Code",IntegerType(),True),
                                StructField("Region Code",IntegerType(),True),
                                StructField("official_name_ar",StringType(),True),
                                StructField("ISO4217-currency_minor_unit",IntegerType(),True),
                                StructField("UNTERM Arabic Formal",StringType(),True),
                                StructField("UNTERM Chinese Short",StringType(),True),
                                StructField("Land Locked Developing Countries
(LLDC)",StringType(),True),
                                StructField("Intermediate Region Name",StringType(),True),
                                StructField("official_name_es",StringType(),True),
                                StructField("UNTERM English Formal",StringType(),True),
                                StructField("official_name_cn",StringType(),True),
                                StructField("official_name_en",StringType(),True),
                                StructField("ISO4217-
currency_country_name",StringType(),True),
                                StructField("Least Developed Countries
(LDC)",StringType(),True),
                                StructField("Region Name",StringType(),True),
                                StructField("UNTERM Arabic Short",StringType(),True),
                                StructField("Sub-region Name",StringType(),True),
                                StructField("official_name_ru",StringType(),True),
                                StructField("Global Name",StringType(),True),
                                StructField("Capital",StringType(),True),
                                StructField("Continent",StringType(),True),
                                StructField("TLD",StringType(),True),
                                StructField("Languages",StringType(),True),
                                StructField("Geoname ID",IntegerType(),True),
                                StructField("CLDR display name",StringType(),True),
                                StructField("EDGAR",StringType(),True),

                    ])
```

```
28          StructField("UNTERM French Format",StringType(),True),
29          StructField("UNTERM Russian Short",StringType(),True),
30          StructField("M49",IntegerType(),True),
31          StructField("Sub-region Code",IntegerType(),True),
32          StructField("Region Code",IntegerType(),True),
33          StructField("official_name_ar",StringType(),True),
34          StructField("ISO4217-currency_minor_unit",IntegerType(),True),
35          StructField("UNTERM Arabic Formal",StringType(),True),
36          StructField("UNTERM Chinese Short",StringType(),True),
37          StructField("Land Locked Developing Countries (LLDC)",StringType(),True),
38          StructField("Intermediate Region Name",StringType(),True),
39          StructField("official_name_es",StringType(),True),
40          StructField("UNTERM English Formal",StringType(),True),
41          StructField("official_name_cn",StringType(),True),
42          StructField("official_name_en",StringType(),True),
43          StructField("ISO4217-currency_country_name",StringType(),True),
44          StructField("Least Developed Countries (LDC)",StringType(),True),
45          StructField("Region Name",StringType(),True),
46          StructField("UNTERM Arabic Short",StringType(),True),
47          StructField("Sub-region Name",StringType(),True),
48          StructField("official_name_ru",StringType(),True),
49          StructField("Global Name",StringType(),True),
50          StructField("Capital",StringType(),True),
51          StructField("Continent",StringType(),True),
52          StructField("TLD",StringType(),True),
53          StructField("Languages",StringType(),True),
54          StructField("Geoname ID",IntegerType(),True),
55          StructField("CLDR display name",StringType(),True),
56          StructField("EDGAR",StringType(),True),
57
58      ])
```

Command took 0.04 seconds -- by dxc262ab1225_1654530016505@manipalazure.onmicrosoft.com at 6/15/2022, 4:45:44 PM on clus1219

```python
country_codes_df = spark.read \
.option("header" , True) \
.schema(country_codes_schema) \
.csv("/FileStore/tables/country_codes.csv")
```



```python
1  country_codes_df = spark.read \
2  .option("header" , True) \
3  .schema(country_codes_schema) \
4  .csv("/FileStore/tables/country_codes.csv")
```

▶ 📇 country_codes_df: pyspark.sql.dataframe.DataFrame = [FIFA: string, Dial: string ... 54 more fields]

Command took 0.28 seconds -- by dxc262ab1225_1654530016505@manipalazure.onmicrosoft.com at 6/15/2022, 4:50:01 PM on clus1219

```python
from pyspark.sql.functions import col,lit

country_codes_selected_df = country_codes_df.select(col('FIFA'),
                                col('Dial'),col('Developed / Developing
Countries').alias('D/UD'),col('UNTERM Chinese
Short').alias('Unterm_Chinese_Short'),col('Land Locked Developing Countries
(LLDC)').alias('LLDC'),col('official_name_es'),col('Region Name'),col('EDGAR'))
```



```python
1  country_codes_selected_df = country_codes_df.select(col('FIFA'),
2                              col('Dial'),col('Developed / Developing Countries').alias('D/UD'),col('UNTERM Chinese
   Short').alias('Unterm_Chinese_Short'),col('Land Locked Developing Countries (LLDC)').alias('LLDC'),col('official_name_es'),col('Region Name'),col('EDGAR'))
```

▼ 📇 country_codes_selected_df: pyspark.sql.dataframe.DataFrame
    FIFA: string
    Dial: string
    D/UD: string
    Unterm_Chinese_Short: string
    LLDC: string
    official_name_es: string
    Region Name: string
    EDGAR: string

Command took 0.08 seconds -- by dxc262ab1225_1654530016505@manipalazure.onmicrosoft.com at 6/15/2022, 5:06:51 PM on clus1219
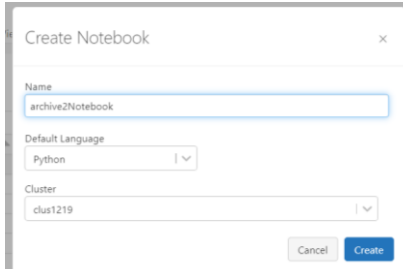
display(country_codes_selected_df)



**Case 2.** Using archive2.zip file - please ingest data into Databricks DBFS path & query the data, redesign columns accordingly using Dataframe commands - display with notebooks accordingly.

**File Being Used: nces330.20.csv**

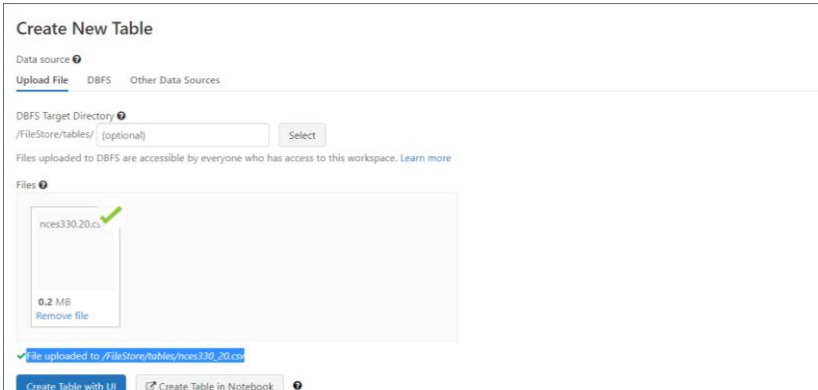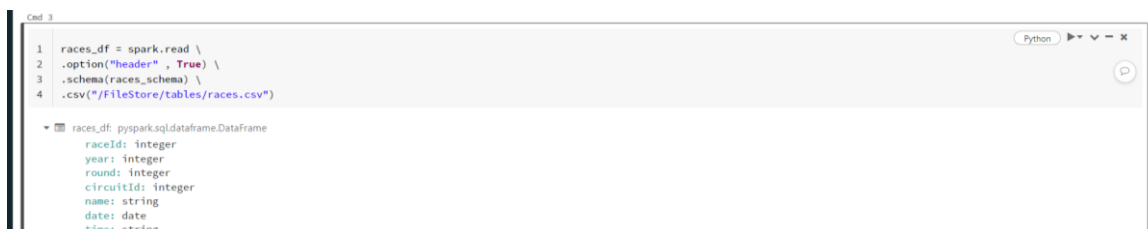**Step 1:** First, login to your Azure Portal and create a Databricks workspace.

**Step 2:** Open the Databricks workspace and create cluster for your future use.

**Step 3:** Now, create a notebook by clicking on the create Notebook option from the side panel.



**Step 4:** After creating the notebook, ingest the data into the Databricks by dragging and dropping the required file in the **drag & drop** region.

```python
from pyspark.sql.types import StructType, StructField, IntegerType, StringType
```

```
Cmd 1
 1  from pyspark.sql.types import StructType, StructField, IntegerType, StringType

Command took 0.04 seconds -- by dxc262ab1225_1654530016505@manipalazure.onmicrosoft.com at 6/15/2022, 5:24:11 PM on clus1219
```

```python
nces330_20_schema = StructType(fields=[StructField("year",IntegerType(),False),
                       StructField("State",StringType(),True),
                       StructField("Type",StringType(),True),
                       StructField("Length",StringType(),True),
                       StructField("Expense",StringType(),True),
                       StructField("Value",IntegerType(),True),
                      ])
```

```
 1  races_schema = StructType(fields=[StructField("raceId",IntegerType(),False),
 2                      StructField("year",IntegerType(),True),
 3                      StructField("round",IntegerType(),True),
 4                      StructField("circuitId",IntegerType(),True),
 5                      StructField("name",StringType(),True),
 6                      StructField("date",DateType(),True),
 7                      StructField("time",StringType(),True),
 8                      StructField("url",StringType(),True),
 9                      ])

Command took 0.04 seconds -- by dxc262ab1225_1654530016505@manipalazure.onmicrosoft.com at 6/15/2022, 3:18:02 PM on clus1219
Cmd 3
```

```python
nces330_20_df = spark.read \
.option("header" , True) \
.schema(nces330_20_schema) \
.csv("/FileStore/tables/nces330_20.csv")
```

```
Cmd 3                                                                    Python  ▶▾ ∨ — ✕
 1  races_df = spark.read \
 2  .option("header" , True) \
 3  .schema(races_schema) \
 4  .csv("/FileStore/tables/races.csv")

 ▾ ▦ races_df: pyspark.sql.dataframe.DataFrame
      raceId: integer
      year: integer
      round: integer
      circuitId: integer
      name: string
      date: date
      time: string
```

```python
from pyspark.sql.functions import col,lit
```

```python
nces330_20_selected_df = nces330_20_df.select(col('Year'),
                             col('State'),col('Expense'))
```

```
 1  from pyspark.sql.functions import col,lit

Command took 0.03 seconds -- by dxc262ab1225_1654530016505@manipalazure.onmicrosoft.com at 6/15/2022, 5:34:02 PM on clus1219
Cmd 5
```
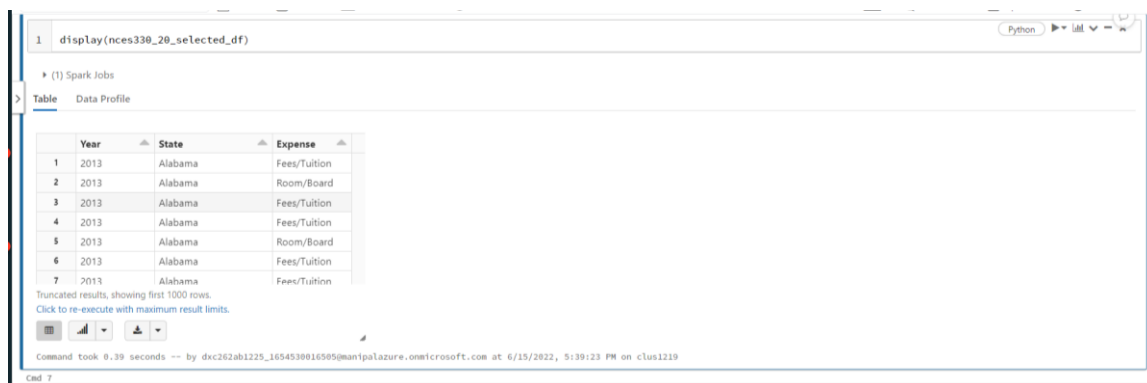
```
 1  nces330_20_selected_df = nces330_20_df.select(col('Year'),          Python  ▶▾ ∨ — ✕
 2                              col('State'),col('Expense'))

 ▾ ▦ nces330_20_selected_df: pyspark.sql.dataframe.DataFrame
      Year: integer
      State: string
      Expense: string

Command took 0.06 seconds -- by dxc262ab1225_1654530016505@manipalazure.onmicrosoft.com at 6/15/2022, 5:37:09 PM on clus1219
Cmd 6
```

display(nces330_20_selected_df)



**Case 3.** Using archive3.zip file - please ingest data into Databricks DBFS path & query the data, redesign columns accordingly using Dataframe commands - display with notebooks accordingly.
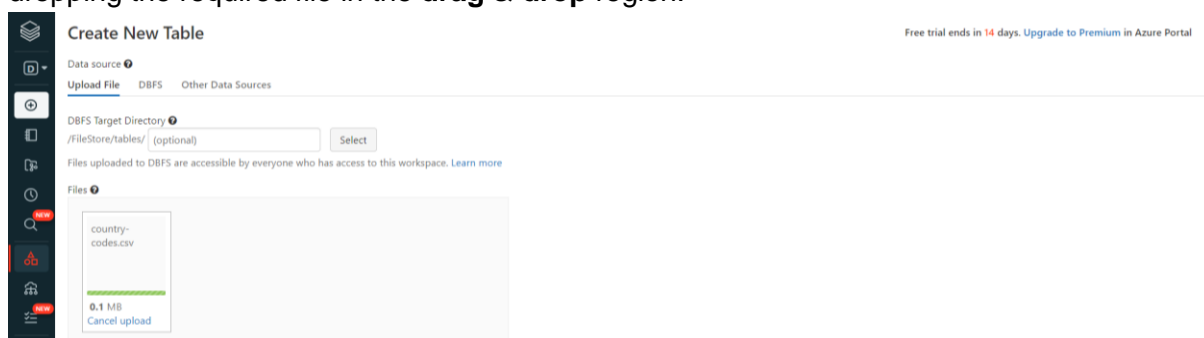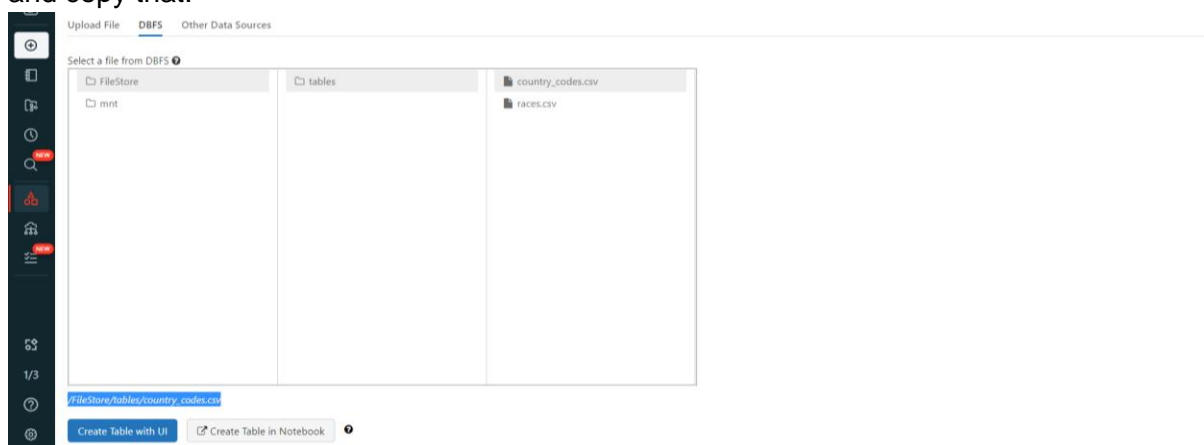
**File Being used: final_data.csv**

**Step 1:** First, login to your Azure Portal and create a Databricks workspace.

**Step 2:** Open the Databricks workspace and create cluster for your future use.

**Step 3:** Now, create a notebook by clicking on the create Notebook option from the side panel.

**Step 4:** After creating the notebook, ingest the data into the Databricks by dragging and dropping the required file in the **drag & drop** region.



Later, click on DBFS and select the file that you have dropped. This will give you the file path and copy that.

**Step-5:** Import the required fields and features from pyspark.

from pyspark.sql.types import StructType, StructField, IntegerType, StringType

final_data_schema = StructType(fields=[StructField("tweet_text",StringType(),False),
                    StructField("emotion_in_tweet_is_directed_at",StringType(),True),

StructField("is_there_an_emotion_directed_at_a_brand_or_product",StringType(),True),
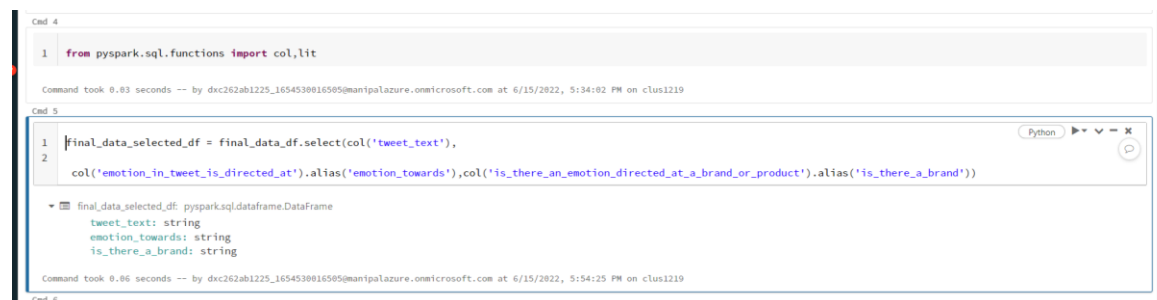                    ])


final_data_df = spark.read \
.option("header" , True) \
.schema(nces330_20_schema) \
.csv("/FileStore/tables/final_data.csv")



from pyspark.sql.functions import col,lit

final_data_selected_df = final_data_df.select(col('tweet_text'),

col('emotion_in_tweet_is_directed_at').alias('emotion_towards'),col('is_there_an_emotion_dir
ected_at_a_brand_or_product').alias('is_there_a_brand'))

display(final_data_selected_df)



**Case 4.** Using archive4.zip file - please ingest data into Databricks DBFS path & query the data, redesign columns accordingly using Dataframe commands - display with notebooks accordingly.
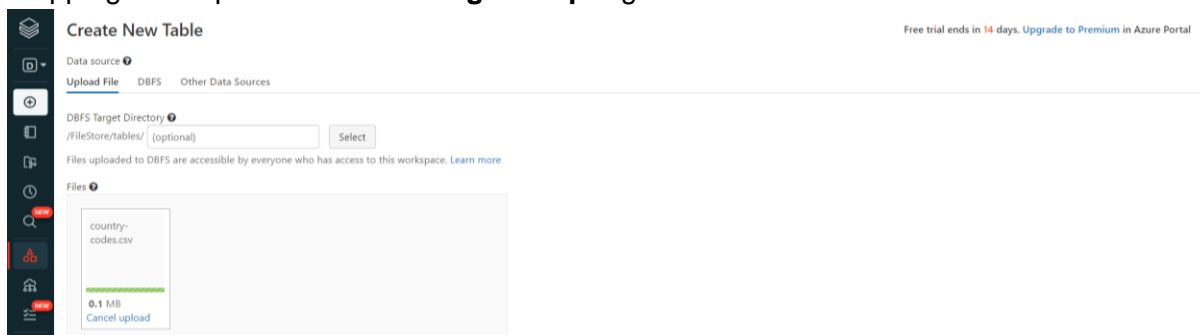
**File Being used: SEntFiN-v1.1.csv**

**Step 1:** First, login to your Azure Portal and create a Databricks workspace.
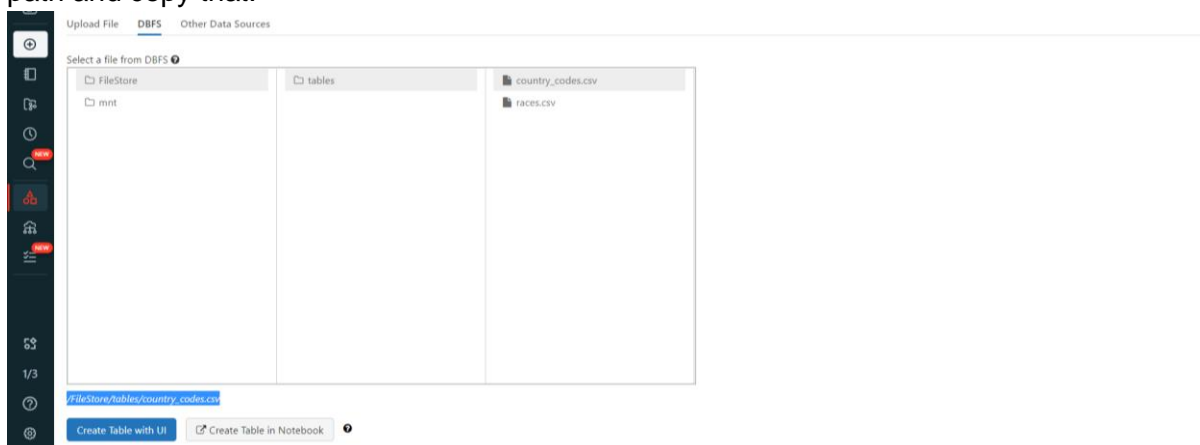
**Step 2:** Open the Databricks workspace and create cluster for your future use.

**Step 3:** Now, create a notebook by clicking on the create Notebook option from the side panel.

**Step 4:** After creating the notebook, ingest the data into the Databricks by dragging and dropping the required file in the **drag & drop** region.



Later, click on DBFS and select the file that you have dropped. This will give you the file path and copy that.

**Step-5:** Import the required fields and features from pyspark.

```
from pyspark.sql.types import StructType, StructField, IntegerType, StringType

from pyspark.sql.types import StructType, StructField, IntegerType, StringType

SEntFiN-v1_1_schema = StructType(fields=[StructField("S No.",IntegerType(),False),
                  StructField("Title",StringType(),True),
                  StructField("Decisions",StringType(),True),
                  StructField("Words",IntegerType(),True),

])


SEntFiN-v1_1_df = spark.read \
.option("header" , True) \
.schema(SEntFiN-v1_1_schema) \
.csv("/FileStore/tables/ SEntFiN-v1_1.csv")

from pyspark.sql.functions import col,lit

SEntFiN-v1_1_selected_df = SEntFiN-v1_1_df.select(col('S No'),
                          col('Title'),col('Words'))

display(SEntFiN-v1_1_df)
```

Case 5. Using archive5.zip file - please ingest data into Databricks DBFS path & query the data, redesign columns accordingly using Dataframe commands - display with notebooks accordingly.
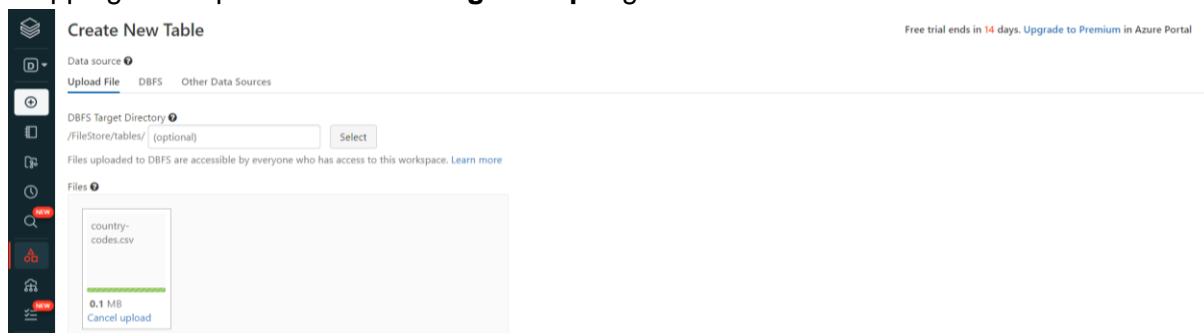
**File Being used: cancer_death_rates.csv**
**Step 1:** First, login to your Azure Portal and create a Databricks workspace.
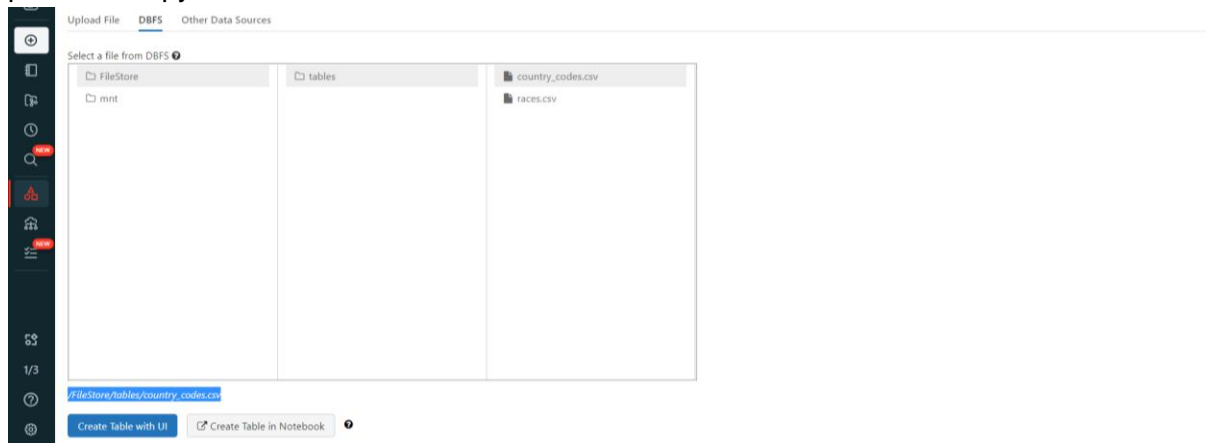**Step 2:** Open the Databricks workspace and create cluster for your future use.
**Step 3:** Now, create a notebook by clicking on the create Notebook option from the side panel.
**Step 4:** After creating the notebook, ingest the data into the Databricks by dragging and dropping the required file in the **drag & drop** region.

Later, click on DBFS and select the file that you have dropped. This will give you the file path and copy that.



**Step-5:** Import the required fields and features from pyspark.

```
from pyspark.sql.types import StructType, StructField, IntegerType, StringType, FloatType

cancer_death_rates_schema = StructType(fields=[StructField("Entity",StringType(),False),
                    StructField("Code",StringType(),True),
                    StructField("Year",IntegerType(),True),
                    StructField("Deaths - Neoplasms - Sex: Both - Age: Age-standardized
(Rate)",FloatType(),True),

])


Cancer_death_rates_df = spark.read \
.option("header" , True) \
.schema(cancer_death_rates_schema) \
.csv("/FileStore/tables/ cancer_death_rates.csv")

from pyspark.sql.functions import col,lit

cancer_death_rates_selected_df = cancer_death_rates_df.select(col(' Entity'),
                        col(' Year'),col(' Deaths - Neoplasms - Sex: Both - Age: Age-
standardized (Rate)').alias('Deaths'))

display(cancer_death_rates_df)
```

Case 6. Using archive6.zip file - please ingest data into Databricks DBFS path & query the data, redesign columns accordingly using Dataframe commands - display with notebooks accordingly.
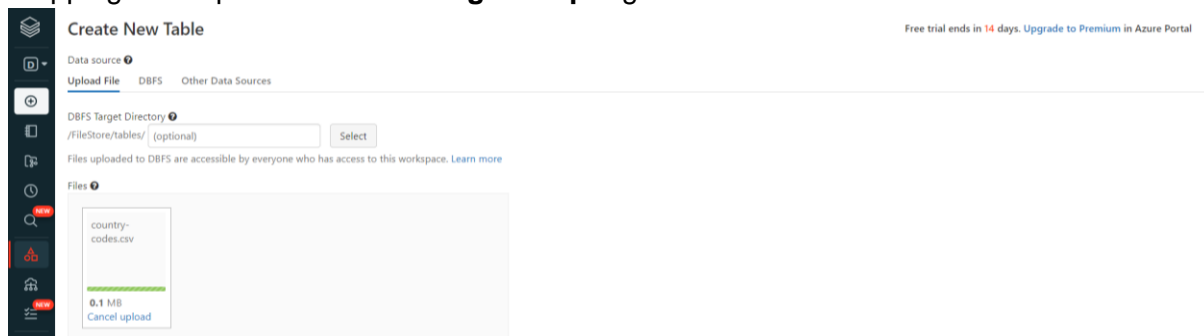
**File Being used: inflation_gdp.csv**
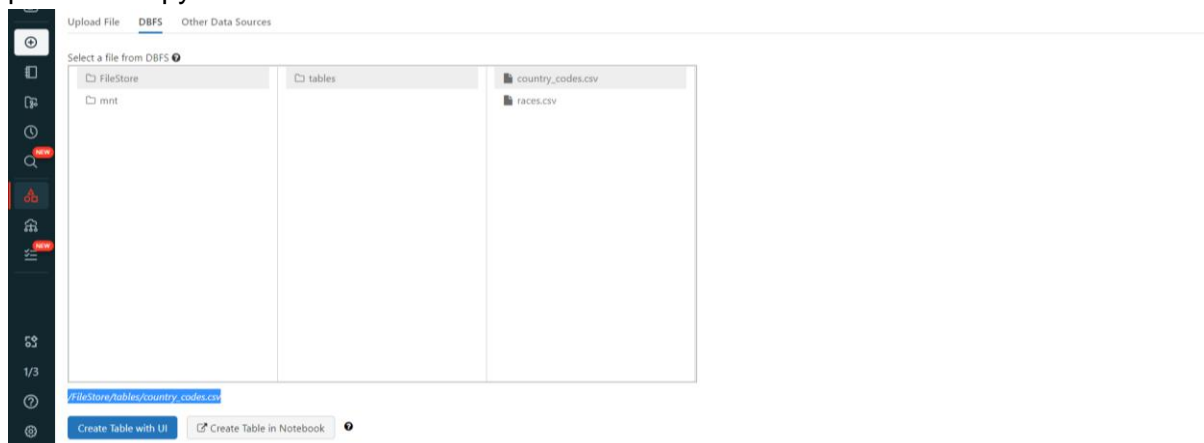**Step 1:** First, login to your Azure Portal and create a Databricks workspace.
**Step 2:** Open the Databricks workspace and create cluster for your future use.
**Step 3:** Now, create a notebook by clicking on the create Notebook option from the side panel.
**Step 4:** After creating the notebook, ingest the data into the Databricks by dragging and dropping the required file in the **drag & drop** region.



 Later, click on DBFS and select the file that you have dropped. This will give you the file path and copy that.



**Step-5:** Import the required fields and features from pyspark.

```
from pyspark.sql.types import StructType, StructField, IntegerType, StringType, FloatType

inflation_gdp_schema = StructType(fields=[StructField("Country ",StringType(),False),
                 StructField("Country Code ",StringType(),True),
                 StructField("Year ",IntegerType(),True),
                 StructField("Inflation ",FloatType(),True),

])


inflation_gdp _df = spark.read \
.option("header" , True) \
```

```
.schema(inflation_gdp _schema) \
.csv("/FileStore/tables/ inflation_gdp.csv")

from pyspark.sql.functions import col,lit

inflation_gdp _selected_df = inflation_gdp _df.select(col('S No'),
                                col('Title'),col('Words'))

display(inflation_gdp _df)
```