

Day25 - March1st 2024

1. Woked up early and started my day
2. Cooked food for mrng and for afternoon
3. Solved one medium leetcode problem of binary search
4. A realtime pyspark project is in progress
5. Faced with some errors solved it on my own

pythonProject1 Youtube Version control

Project

- pythonProject1 Youtube
  - Properties
  - configuration
    - logging.config
  - source
  - olap
  - oltp
  - venv library root
    - create\_spark.py
    - driver.py
    - get\_env\_variables.py
    - main.py
    - validate.py
  - python C:\spark\1\spark-3.3.3-bin-hadoop3.zip
  - External Libraries
  - Scratches and Consoles

Run driver

Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties  
Setting default log level to "WARN".  
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel.  
object created... <pyspark.sql.session.SparkSession object at 0x00002320A80190>  
validating spark object with current date=[Row(current\_date=datetime.date(2023, 4, 15))]  
Process finished with exit code 0

Implementing Pyspark Real Time Application || End-to-End

DataSpark  
1.32K subscribers

Subscribed

436

13K views 8 months ago PYSPARK REAL TIME APPLICATION  
In this video we will discuss about , implementing Pyspark application in Pycharm Dynamically from the Respective Folders..

4°C Partly sunny

pythonProject1 Version control

Project

- pythonProject1
  - venv library root
    - Lib
    - Scripts
    - share
    - gitignore
    - pyvenv.cfg
  - Properties
    - configuration
      - logging.config
  - source
    - olap
    - oltp
    - create\_spark.py
    - driver.py
    - get\_env\_variables.py
    - main.py
    - validate.py
  - python C:\spark\1\spark-3.3.3-bin-hadoop3.zip

Run driver

pythonProject1\venv\scripts\python.exe -u pythonProject1\driver.py  
I am in main  
24/03/01 09:37:39 WARN Shell: Did not find winutils.exe: java.io.FileNotFoundException  
Setting default log level to "WARN".  
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel.  
object created... <pyspark.sql.session.SparkSession object at 0x000027A0870B7A0>  
validating spark object with current date=[Row(current\_date=datetime.date(2023, 4, 15))]  
Process finished with exit code 0

6.

pythonProject1 Version control

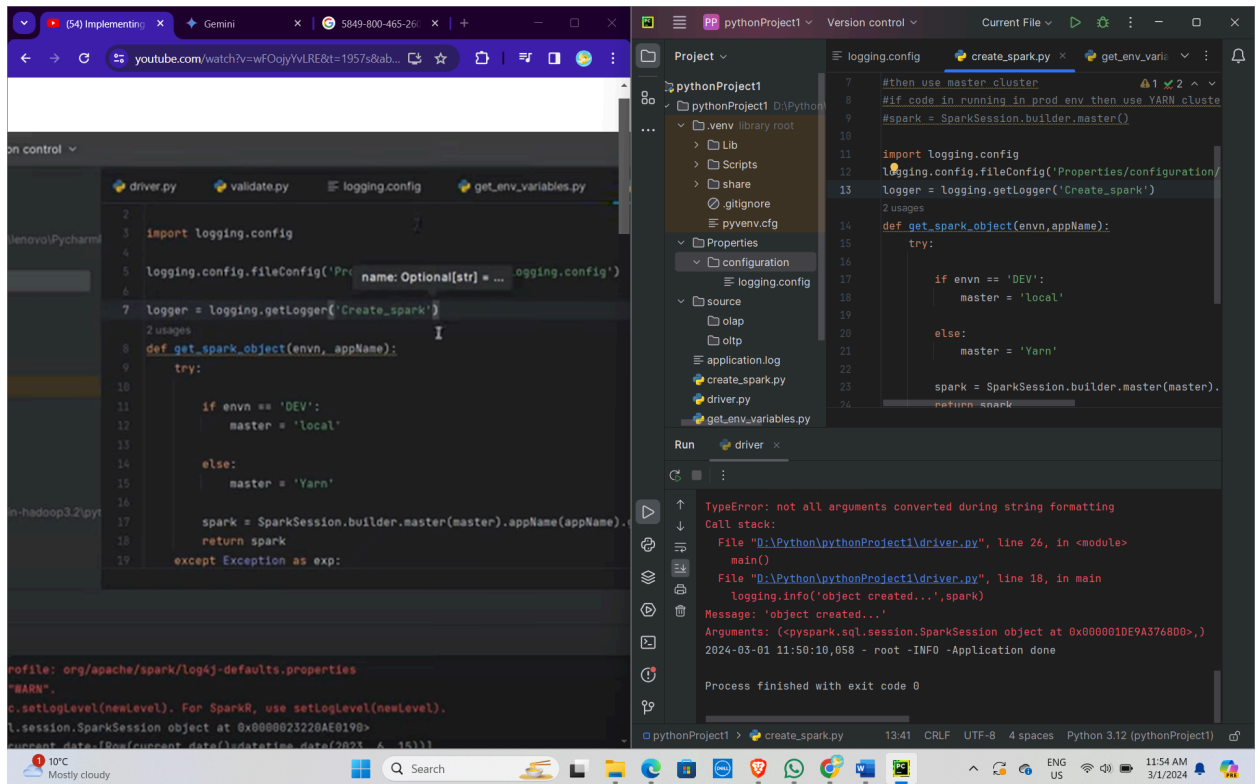
Project

- pythonProject1
  - venv library root
    - Lib
    - Scripts
    - share
    - gitignore
    - pyvenv.cfg
  - Properties
    - configuration
      - logging.config
  - source
    - olap
    - oltp
    - application.log
    - create\_spark.py
    - driver.py
    - get\_env\_variables.py

Run driver

TypeError: not all arguments converted during string formatting  
Call stack:  
File "D:\Python\pythonProject1\driver.py", line 26, in <module>  
main()  
File "D:\Python\pythonProject1\driver.py", line 18, in main  
logging.info('object created...',spark)  
Message: 'object created...'  
Arguments: (<pyspark.sql.session.SparkSession object at 0x000001DE9A376BD0>,)  
2024-03-01 11:50:10,058 - root -INFO -Application done  
Process finished with exit code 0

```
1 import logging.config
2
3 logging.config.fileConfig('Properties/configuration/logging.config')
4
5
6
7
8 def get_spark_object(envn, appName):
9     try:
10
11         if envn == 'DEV':
12             master = 'local'
13
14         else:
15             master = 'Yarn'
16
17         spark = SparkSession.builder.master(master).appName(appName).getOrCreate()
18         return spark
19     except Exception as exp:
```



7. Ended my day by solving a real world SQL problem from Ankit's YT

SQLQuery1.sql - KAUSHI\SQLEXPRESS.master (KAUSHI\jamka (62)) - Microsoft SQL Server Management Studio

File Edit View Query Project Tools Window Help

master

Execute

```
select * from event_status;

with cte as(
select event_time,status+lag(status,1,status) over(order by event_time) as l from event_status)
,cte2 as(
select *,sum(case when l = 'onoff' then 1 else 0 end) over(order by event_time) as rn
from cte)
select min(event_time) as login,max(event_time) as logout,count(1)-1 as onCount
from cte2
group by rn

/* Explanation:
Step1 : Get previous status using lag
Step2 : In order to create an id for the sequence of on/off, use case within SUM to create a group key.
This is basically continuously checking if the order is changing from previous value.
Step3 : Use the group key(rn) to get anything.
```

109 %

Results Messages

	event_time	status
1	10:01	on
2	10:02	on
3	10:03	on
4	10:04	off
5	10:07	on
6	10:08	on
7	10:09	off
8	10:11	on
9	10:12	off

	login	logout	onCount
1	10:01	10:04	3
2	10:07	10:09	2
3	10:11	10:12	1

Query executed successfully.

KAUSHI\SQLEXPRESS (16.0 RTM) KAUSHI\jamka (62) master 00:00:00 12 rows

Ready

10°C Clear

Ln 27 Col 11 Ch 11 INS

7:40 PM 3/1/2024

SQLQuery1.sql - KAUSHI\SQLEXPRESS.master (KAUSHI\jamka (62)) - Microsoft SQL Server Management Studio

File Edit View Query Project Tools Window Help

master

Execute

```
select * from event_status;

with cte as(
select event_time,status+lag(status,1,status) over(order by event_time) as l from event_status)
,cte2 as(
select *,sum(case when l = 'onoff' then 1 else 0 end) over(order by event_time) as rn
from cte)
select min(event_time) as login,max(event_time) as logout,count(1)-1 as onCount
from cte2
group by rn

/* Explanation:
Step1 : Get previous status using lag
Step2 : In order to create an id for the sequence of on/off, use case within SUM to create a group key.
This is basically continuously checking if the order is changing from previous value.
Step3 : Use the group key(rn) to get anything.
```

109 %

Messages

Msg 113, Level 15, State 1, Line 35  
Missing end comment mark '\*/'.

Completion time: 2024-03-01T22:50:03.0679867-06:00

Query completed with errors.

KAUSHI\SQLEXPRESS (16.0 RTM) KAUSHI\jamka (62) master 00:00:00 0 rows

Ready

7°C Clear

Ln 29 Col 1 Ch 1 INS

10:50 PM 3/1/2024

SQLQuery1.sql - KAUSHI\SQLEXPRESS.master (KAUSHI\janka (62)) - Microsoft SQL Server Management Studio

```

select * from event_status;

with cte as(
select event_time,status,lag(status,1,status) over(order by event_time) as l from event_status
),cte2 as(
select *,sum(case when l = 'onoff' then 1 else 0 end) over(order by event_time) as rn
from cte)
select min(event_time) as login,max(event_time) as logout,count(1) as onCount
from cte2
group by rn

/* Explanation:
Step1 : Get previous status using lag
Step2 : In order to create an id for the sequence of on/off, use case within SUM to create a group key.
This is basically continuously checking if the order is changing from previous value.
Step3 : Use the group key(rn) to get anything.

```

Results

event_time	status
10:01	on
10:02	on
10:03	on
10:04	off
10:07	on
10:08	on
10:09	off
10:11	on
10:12	off

login	logout	onCount
10:01	10:04	3
10:07	10:09	2
10:11	10:12	1

Query executed successfully.

SQLQuery3.sql - KAUSHI\SQLEXPRESS.master (KAUSHI\janka (57)) - Microsoft SQL Server Management Studio

```

insert into players_location
values ('Sachin','Mumbai'),('Virat','Delhi'),('Rahul','Bangalore'),('Rohit','Mumbai'),('Mayank','Bangalore');

select * from players_location;

with cte as(
select *,row_number() over(partition by city order by name) as rn
from players_location)
select min(case when city = 'Bangalore' then name end) as Bangalore,
min(case when city = 'Mumbai' then name end) as Mumbai,
min(case when city = 'Delhi' then name end) as Delhi
from cte
group by rn

/*Explanation:
Step1 : First we have used row_number over by city order by name
Step2 : Now using rn as key we will perform group by operation
Step3 : Then we use case when and aggregate functions to get the output as required

```

Results

name	city
Sachin	Mumbai
Virat	Delhi
Rahul	Bangalore
Rohit	Mumbai
Mayank	Bangalore

Bangalore	Mumbai	Delhi
Mayank	Rohit	Virat
Rahul	Sachin	NULL

Query executed successfully.