

# ADVANCED SPARK SQL TECHNIQUES: EXPLORING WINDOW CLAUSE

Unleashing  
the Power of  
Data  
Analytics



# CHALLENGE AT HAND : SEQUENTIAL AGG

- WINDOW FUNCTION IS USEFUL WHEN THE CHANGE ACROSS TIME NEEDS TO BE ANALYSED
- RANKING THE ENTITIES RELATIVE TO TIME AND ARBITRARY DATA POINT
- UNLIKE GROUP BY NEW COLUMN IS CREATED FROM WINDOW OPERATION
- 3 DIFFERENT FUNCTIONS:

## RANKING FUNCTIONS

- SYNTAX: RANK | DENSE\_RANK | PERCENT\_RANK | NTILE | ROW\_NUMBER

## ANALYTIC FUNCTIONS

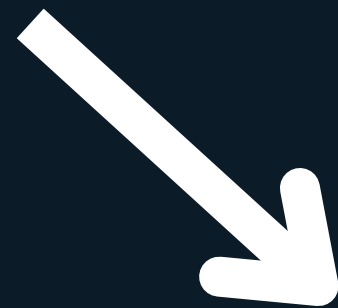
- SYNTAX: CUME\_DIST | LAG | LEAD | NTH\_VALUE | FIRST\_VALUE | LAST\_VALUE

## AGGREGATE FUNCTIONS

- SYNTAX: MAX | MIN | COUNT | SUM | AVG | ...

# HOW WE ARE DOING IT?

USE KAGGLE  
NOTEBOOK TO  
LOAD DATA IN  
PYSPARK



ABOVE  
COMMANDS  
ARE EXECUTED

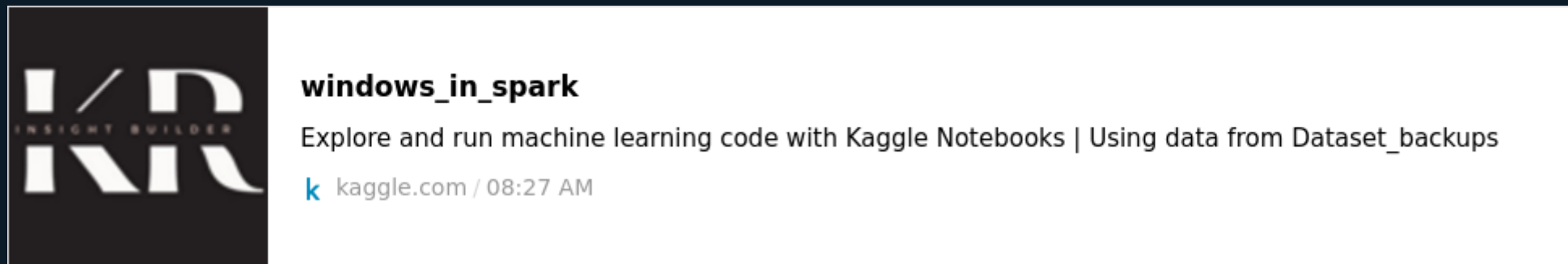


DISCUSS THE  
RESULTS OF  
EXECUTION AND  
PROBLEM SOLVED



TROUBLE SHOOTING  
ISSUES THAT ARISES  
IN MESSY DATA

# LETS GET OURSELF A PYSPARK NOTEBOOK AND DIG IN



REAL CLUSTER IS NOT  
NECESSARY FOR LEARNING  
THE DML

# THANKS FOR WATCHING

**PRACTICE**

**PRACTICE**

 **LIKE**

 **SHARE**

 **SUBSCRIBE**

**PRACTICE**

**PRACTICE**