

DataFrame ด้วย PySpark และ Pandas

```
import pandas as pd
from pyspark.sql import SparkSession
```

```
# สร้าง Pandas DataFrame
```

```
pandas_df = pd.DataFrame({
    "name": ["Alice", "Bob", "Charlie"],
    "age": [25, 30, 35]
})
```

```
# แปลง Pandas DataFrame ไปเป็น PySpark DataFrame
spark_df = spark.createDataFrame(pandas_df)
spark_df.show()
```

```
df.show(5) # แสดง 5 แถวแรก
```

```
df.printSchema() # แสดงโครงสร้างของ DataFrame
```

```
df.select("name", "age").show()
```

```
df.filter(df.age > 30).show() # แสดงข้อมูลที่มีอายุมากกว่า 30
```

```
df.groupBy("age").count().show() # นับจำนวนของแต่ละกลุ่มอายุ
```

```
from pyspark.sql.functions import col
```

```
# เพิ่มคอลัมน์ใหม่ที่ชื่อว่า 'age_plus_10' ซึ่งเพิ่มอายุขึ้น 10 ปี
df = df.withColumn("age_plus_10", col("age") + 10)
df.show()
```

```
+-----+---+-----+
|  name|age|age_plus_10|
+-----+---+-----+
|  Alice| 25|         35|
|   Bob| 30|         40|
|Charlie| 35|         45|
+-----+---+-----+
```

```

from pyspark.sql.functions import when

# เพิ่มคอลัมน์ 'age_group' เพื่อกำหนดกลุ่มอายุ
df = df.withColumn("age_group", when(col("age") < 30, "Young")
                                   .when((col("age") >= 30) & (col("age") < 40), "Middle-aged")
                                   .otherwise("Old"))

df.show()

```

```

+-----+---+-----+-----+
|  name|age|age_plus_10|  age_group|
+-----+---+-----+-----+
| Alice| 25|         35|    Young|
|  Bob| 30|         40|Middle-aged|
|Charlie| 35|         45|Middle-aged|
+-----+---+-----+-----+

```

```

from pyspark.sql.functions import avg, count, max, min

# หาอายุเฉลี่ย
df_avg = df.agg(avg("age").alias("average_age"))
df_avg.show()

# นับจำนวนคนในแต่ละกลุ่มอายุ
df_group_count = df.groupBy("age_group").agg(count("*").alias("count"))
df_group_count.show()

```

```

+-----+
|average_age|
+-----+
|         30.0|
+-----+

```

```

# บันทึกเป็น CSV
df.write.csv("path/to/output.csv", header=True)

# บันทึกเป็น Parquet
df.write.parquet("path/to/output.parquet")

```

Graph Analytics PowerBI

```
1 from pyspark.sql import SparkSession
2 from graphframes import GraphFrame
3 from pyspark.sql.functions import desc, col, lit
4
5 spark = SparkSession.builder \
6     .appName("CyclingRoutesGraph") \
7     .config("spark.jars.packages", "graphframes:graphframes:0.8.2-spark3.0-s_2.12") \
8     .getOrCreate()
9
10 cycling_routes_df = spark.read.csv("cycling.csv", header=True, inferSchema=True)
11
12 vertices = cycling_routes_df.select("FromStationName").withColumnRenamed("FromStationName", "id").distinct()
13
14 edges = cycling_routes_df.select("FromStationName", "ToStationName") \
15     .withColumnRenamed("FromStationName", "src") \
16     .withColumnRenamed("ToStationName", "dst")
17
18 graph = GraphFrame(vertices, edges)
19
20 grouped_edges = graph.edges.groupBy("src", "dst").count() \
21     .filter(col("count") >= 3) \
22     .withColumn("source_color", lit("#FF3F33")) \
23     .withColumn("destination_color", lit("#3358FF"))
24
25 grouped_edges.select("src", "dst", "source_color", "destination_color") \
26     .write.csv("grouped_cycling_routes_no_header.csv", mode="overwrite", header=False)
27
28 grouped_edges.show()
29
30 print("Processing and writing to file completed.")
31
```

```
warnings.warn(
+-----+-----+-----+-----+-----+
+|          src|          dst|count|source_color|destination_color|
+-----+-----+-----+-----+-----+
| Union St & 4th Ave| 3rd Ave & Broad St| 8| #FF3F33| #3358FF|
| 2nd Ave & Pine St|Occidental Park /...| 34| #FF3F33| #3358FF|
| NE 42nd St & Univ...|UW Intramural Act...| 6| #FF3F33| #3358FF|
| 15th Ave E & E Th...|NE Pacific St/UW ...| 3| #FF3F33| #3358FF|
| Dexter Ave N & Al...|6th Ave & Blancha...| 8| #FF3F33| #3358FF|
| Pier 69 / Alaskan...| 3rd Ave & Broad St| 9| #FF3F33| #3358FF|
| 15th Ave E & E Th...|Thomas St & 5th A...| 3| #FF3F33| #3358FF|
| REI / Yale Ave N ...|Dexter Ave N & Al...| 15| #FF3F33| #3358FF|
| 12th Ave & E Denn...|Occidental Park /...| 4| #FF3F33| #3358FF|
| E Blaine St & Fai...|15th Ave NE & NE ...| 3| #FF3F33| #3358FF|
| 15th Ave NE & NE ...|Eastlake Ave E & ...| 7| #FF3F33| #3358FF|
| 2nd Ave & Vine St|Key Arena / 1st A...| 12| #FF3F33| #3358FF|
| E Pine St & 16th Ave|Bellevue Ave & E ...| 15| #FF3F33| #3358FF|
| PATH / 9th Ave & ...| Union St & 4th Ave| 16| #FF3F33| #3358FF|
| King Street Stati...| Pine St & 9th Ave| 6| #FF3F33| #3358FF|
| Thomas St & 5th A...| Union St & 4th Ave| 4| #FF3F33| #3358FF|
| Summit Ave & E De...| Union St & 4th Ave| 6| #FF3F33| #3358FF|
| E Harrison St & B...|Broadway and E De...| 13| #FF3F33| #3358FF|
| Cal Anderson Park...|Broadway and E De...| 6| #FF3F33| #3358FF|
| 2nd Ave & Blancha...|Thomas St & 5th A...| 9| #FF3F33| #3358FF|
+-----+-----+-----+-----+-----+
only showing top 20 rows
```

Output CSV

A	B	C	D
Union St & 3rd Ave & E	#FF3F33	#3358FF	
2nd Ave & Occidenta	#FF3F33	#3358FF	
NE 42nd St UW Intram	#FF3F33	#3358FF	
15th Ave E NE Pacific	#FF3F33	#3358FF	
Dexter Ave 6th Ave & E	#FF3F33	#3358FF	
Pier 69 / Al 3rd Ave & E	#FF3F33	#3358FF	
15th Ave E Thomas St	#FF3F33	#3358FF	
REI / Yale A Dexter Ave	#FF3F33	#3358FF	
12th Ave & Occidenta	#FF3F33	#3358FF	
E Blaine St 15th Ave N	#FF3F33	#3358FF	
15th Ave N Eastlake A	#FF3F33	#3358FF	
2nd Ave & Key Arena	#FF3F33	#3358FF	
E Pine St & Bellevue A	#FF3F33	#3358FF	
PATH / 9th Union St &	#FF3F33	#3358FF	

Graph Analytics

```
1 import os
2 os.environ["JAVA_HOME"] = r"C:\Program Files\Java\jdk-22"
3 os.environ["SPARK_HOME"] = r"C:\spark\spark-3.5.1-bin-hadoop3"
4
5 from pyspark.sql import SparkSession
6 from graphframes import GraphFrame
7 from pyspark.sql.functions import desc
8
9 # สร้าง SparkSession สำหรับใช้งาน PySpark
10 spark = SparkSession.builder \
11     .appName("Graph Analytics") \
12     .config("spark.jars.packages", "graphframes:graphframes:0.8.2-spark3.0-s_2.12") \
13     .config("spark.executor.memory", "4g") \
14     .config("spark.driver.memory", "4g") \
15     .config("spark.pyspark.python", "python") \
16     .config("spark.pyspark.driver.python", "python") \
17     .getOrCreate() # สร้างหรือดึง SparkSession
18
19 # ข้อมูลเวอร์เท็กซ์ (จุดยอด) ที่มีชื่อและอายุ
20 vertices_data = [
21     ("Alice", 45),
22     ("Jacob", 43),
23     ("Roy", 21),
24     ("Ryan", 49),
25     ("Emily", 24),
26     ("Sheldon", 52)
27 ]
28 # ข้อมูลเอดจ์ (เชื่อมโยง) ที่มีความสัมพันธ์ระหว่างเวอร์เท็กซ์
29 edges_data = [
30     ("Sheldon", "Alice", "Sister"),
31     ("Alice", "Jacob", "Husband"),
32     ("Emily", "Jacob", "Father"),
33     ("Ryan", "Alice", "Friend"),
34     ("Alice", "Emily", "Daughter"),
35     ("Alice", "Roy", "Son"),
36     ("Jacob", "Roy", "Son")
37 ]
38
39 # สร้าง DataFrame สำหรับเวอร์เท็กซ์
40 vertices_df = spark.createDataFrame(vertices_data, ["id", "age"])
41 # สร้าง DataFrame สำหรับเอดจ์
42 edges_df = spark.createDataFrame(edges_data, ["src", "dst", "relationship"])
43
44 # สร้างกราฟจาก DataFrame ของเวอร์เท็กซ์และเอดจ์
45 graph = GraphFrame(vertices_df, edges_df)
46
47 # แสดงเอดจ์ที่จัดกลุ่มและเรียงตามจำนวน
48 print("Grouped and Ordered Edges:")
49 graph.edges.groupBy("src", "dst").count().orderBy(desc("count")).show(truncate=False)
50
51 # แสดงเอดจ์ที่เกี่ยวข้องกับ 'Alice'
52 print("Filtered Edges where src or dst is 'Alice':")
53 graph.edges.where("src = 'Alice' OR dst = 'Alice'").groupBy("src", "dst").count().orderBy(desc("count")).show(truncate=False)
54
55 # สร้างซับกราฟที่มี 'Alice' เป็นส่วนหนึ่ง
56 print("Subgraph where 'Alice' is involved:")
57 subgraph_edges = graph.edges.where("src = 'Alice' OR dst = 'Alice'")
58 subgraph = GraphFrame(graph.vertices, subgraph_edges)
59 subgraph.edges.show(truncate=False)
60
61 # แสดง motifs ในกราฟที่เกี่ยวข้องกับ 'Alice'
62 print("Motifs in the Graph (connections involving Alice):")
63 motifs = graph.find("(a)-[ab]->(b)")
64 motifs_filtered = motifs.filter("ab.relationship = 'Friend' OR ab.relationship = 'Daughter'")
65 motifs_filtered.show(truncate=False)
66
```

```

67 # คำนวณ PageRank
68 print("PageRank Results:")
69 page_rank = graph.pageRank(resetProbability=0.15, maxIter=5)
70 page_rank.vertices.orderBy(desc("pagerank")).show(truncate=False)
71
72 # แสดง In-Degree ของแต่ละเวอร์เท็กซ์
73 print("In-Degree of Each Vertex:")
74 in_degree = graph.inDegrees
75 in_degree.orderBy(desc("inDegree")).show(truncate=False)
76
77 # ตั้งค่าพ้อยเช็คพอยสำหรับการประมวลผล
78 spark.sparkContext.setCheckpointDir("/tmp/checkpoints")
79 # แสดง Connected Components
80 print("Connected Components:")
81 connected_components = graph.connectedComponents()
82 connected_components.show(truncate=False)
83
84 # แสดง Strongly Connected Components
85 print("Strongly Connected Components:")
86 strongly_connected_components = graph.stronglyConnectedComponents(maxIter=5)
87 strongly_connected_components.show(truncate=False)
88
89 # ค้นหาเส้นทางจาก 'Alice' ไป 'Roy' โดยใช้ BFS
90 print("Breadth-First Search (BFS):")
91 bfs_results = graph.bfs(fromExpr="id = 'Alice'", toExpr="id = 'Roy'", maxPathLength=2)
92 bfs_results.show(truncate=False)
93
94 # สรุปว่าโค้ดนี้ทำการวิเคราะห์กราฟด้วย GraphFrames โดยเริ่มจากการสร้างกราฟจากข้อมูลเวอร์เท็กซ์และเอดจ์ แล้วดำเนินการวิเคราะห์ต่างๆ
95 # เช่น การค้นหาเอดจ์ที่เกี่ยวข้องกับ 'Alice', คำนวณ PageRank, และหาส่วนประกอบที่เชื่อมโยงกันในกราฟ
96 print('END')
97
98 # ปิด Spark session
99 spark.stop()
100

```

Grouped and Ordered Edges:

src	dst	count
Alice	Jacob	1
Sheldon	Alice	1
Emily	Jacob	1
Alice	Emily	1
Alice	Roy	1
Jacob	Roy	1
Ryan	Alice	1

Filtered Edges where src or dst is 'Alice':

src	dst	count
Alice	Jacob	1
Sheldon	Alice	1
Alice	Emily	1
Alice	Roy	1
Ryan	Alice	1

Subgraph where 'Alice' is involved:

src	dst	relationship
Sheldon	Alice	Sister
Alice	Jacob	Husband
Ryan	Alice	Friend
Alice	Emily	Daughter
Alice	Roy	Son

Motifs in the Graph (connections involving Alice):

/usr/local/lib/python3.10/dist-packages/pyspark/sql/dataframe.py:147: UserWarning: DataFrame constructor is internal. Do not directly use it.

a	ab	b
{Ryan, 49}	{Ryan, Alice, Friend}	{Alice, 45}
{Alice, 45}	{Alice, Emily, Daughter}	{Emily, 24}

PageRank Results:

id	age	pagerank
Roy	21	1.9089989375092518
Jacob	43	1.3728466605994618
Alice	45	1.135192093597289
Emily	24	0.7420792759997091
Sheldon	52	0.42044151614714403
Ryan	49	0.42044151614714403

In-Degree of Each Vertex:

id	inDegree
Jacob	2
Alice	2
Roy	2
Emily	1

Connected Components:

id	age	component
Alice	45	197568495616
Jacob	43	197568495616
Roy	21	197568495616
Ryan	49	197568495616
Emily	24	197568495616
Sheldon	52	197568495616

Strongly Connected Components

id	age	component
Roy	21	197568495616
Jacob	43	395136991232
Sheldon	52	730144440320
Emily	24	1022202216448
Alice	45	1065151889408
Ryan	49	1477468749824

Breadth-First Search (BFS):

from	e0	to
{Alice, 45}	{Alice, Roy, Son}	{Roy, 21}

คำสั่ง `groupBy("src", "dst").count()` ใช้ในการนับจำนวนการเชื่อมต่อระหว่างสถานีบนดินทาง และปลายทาง

มีการใช้คำสั่ง `.filter(col("count") > 5)` เพื่อกรองเส้นทางที่มีการบินเข้ามากกว่า 5 ครั้ง

คำสั่ง `.orderBy(desc("count"))` ทำการเรียงลำดับจากมากไปน้อย

ส่วนสี (source_color, destination_color) ถูกเพิ่มเข้าไปใน DataFrame เพื่อใช้ในงานแสดงผล เช่น การวิเคราะห์กราฟขึ้นสูง หรือการ visualizing กราฟ.

ASC เรียงจากน้อยไปหามาก

DESC เรียงจากมากไปหาน้อย

`.orderBy` เรียง

`.filter(col("distance") > 10)` # ตัวอย่างกรองระยะทางมากกว่า 10 กิโลเมตร

`.filter(col("trip_date").between('2023-01-01', '2023-12-31'))` # ตัวอย่างกรองการเดินทางในปี 2023

`.filter(col("src") != col("dst"))` # ตัวอย่างกรองเส้นทางที่มีจุดเริ่มต้นและปลายทางไม่ซ้ำกัน

`hub_stations = graph.edges.groupBy("src").count().orderBy(desc("count"))`

ตัวอย่างหาสถานีที่มีการออกเดินทางมากที่สุด

`.filter(col("trip_type") == "leisure")` # ตัวอย่างกรองเฉพาะการเดินทางเพื่อพักผ่อน

```
grouped_edges = graph.edges.groupBy("src", "dst").count()
.filter(col("count") > 5)
.filter(col("src") != col("dst")) # กรองเฉพาะเส้นทางที่ไม่มีการย้อนกลับ
.filter(col("distance") > 10) # กรองเฉพาะเส้นทางที่มีระยะทางมากกว่า 10 กิโลเมตร
.orderBy(desc("count"))
.withColumn("source_color", lit("#3358FF"))
.withColumn("destination_color", lit("#FF3F33"))
```

หาสถานีที่มีการเดินทางออกจากมากที่สุด

```
popular_start_stations = graph.edges.groupBy("src").count().orderBy(desc("count"))
```

หาสถานีที่มีการเดินทางเข้ามามากที่สุด

```
popular_end_stations = graph.edges.groupBy("dst").count().orderBy(desc("count"))
```

```
popular_start_stations.show()
```

```
popular_end_stations.show()
```

```
popular_routes = graph.edges.groupBy("src","dst").count().orderBy(desc("count"))
```

หาสถานีที่มีเส้นทางเชื่อมโยงกับหลายสถานีที่สุด

```
station_connectivity = graph.degrees.orderBy(desc("degree"))
```

การ Join Table แบบ Pandas Output เอาเป็นไฟล์ csv

```
1 import pandas as pd
2
3 # Load the two uploaded CSV files
4 file1_path = 'fb_live_thailand2.csv'
5 file2_path = 'fb_live_thailand3.csv'
6
7 df1 = pd.read_csv(file1_path)
8 df2 = pd.read_csv(file2_path)
9
10 # Performing an outer join on the dataframes
11 # Assuming default join on common columns
12 outer_joined_df = pd.merge(df1, df2, how='outer') # how='outer,inner,left,right'
13
14 # Display the joined data to the user
15 print(outer_joined_df)
16
17 output_path = 'outer_joined_fb_live_thailand.csv' # สร้างไฟล์ใหม่
18 outer_joined_df.to_csv(output_path, index=False)
19
```


ผลลัพธ์ index=False

	A	B	C	D	E	F	G
1	Table	status_id	status_type	status_public	num_react	num_comments	num_shares
2	FB2	246675545	photo	#####	145	9	
3	FB2	246675545	video	#####	90	78	
4	FB2	246675545	video	#####	75	36	
5	FB2	246675545	video	#####	37	0	
6	FB2	246675545	photo	#####	102	0	
7	FB2	246675545	video	#####	18	0	
8	FB2	246675545	photo	#####	98	0	
9	FB2	246675545	photo	#####	227	7	
10	FB2	246675545	photo	#####	234	15	
11	FB2	246675545	photo	#####	152	2	
12	FB2	246675545	photo	#####	203	1	

ผลลัพธ์ index=True

	A	B	C	D	E	F	G	H
1		Table	status_id	status_type	status_public	num_react	num_comments	num_shares
2	0	FB2	246675545	photo	#####	145	9	
3	1	FB2	246675545	video	#####	90	78	
4	2	FB2	246675545	video	#####	75	36	
5	3	FB2	246675545	video	#####	37	0	
6	4	FB2	246675545	photo	#####	102	0	
7	5	FB2	246675545	video	#####	18	0	
8	6	FB2	246675545	photo	#####	98	0	
9	7	FB2	246675545	photo	#####	227	7	
10	8	FB2	246675545	photo	#####	234	15	
11	9	FB2	246675545	photo	#####	152	2	

การ Join Table แบบ PySpark ไม่มี Output csv

```

1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import trim, col # ใช้ในการตัดช่องว่างและเลือกคอลัมน์
3 from pyspark.sql.types import IntegerType # ใช้สำหรับแปลงข้อมูลให้เป็นชนิด Integer
4
5 # สร้าง SparkSession
6 spark = SparkSession.builder.appName("OuterJoinFillNA").getOrCreate()
7
8 # โหลดข้อมูลจากไฟล์ CSV
9 file_path_1 = "fb_live_thailand2.csv"
10 file_path_2 = "fb_live_thailand3.csv"
11
12 df1 = spark.read.csv(file_path_1, header=True, inferSchema=True)
13 df2 = spark.read.csv(file_path_2, header=True, inferSchema=True)
14
15 join_key = "status_id" # แทนที่ด้วยชื่อคอลัมน์ที่ใช้เป็น key
16
17 data = df1.join(df2, on=join_key, how="outer") # inner , left , right , outer
18
19 print(data)
20
21 fill_data = data.fillna({"num_comments": 0}, {"num_shares": 0})
22
23 # แสดงข้อมูลที่เติมค่า null เป็น 0
24 fill_data.show()
25
26

```


PySpark จัดการ Datframe โดยใช้ภาษา SQL

```
# Import necessary libraries
from pyspark.sql import SparkSession

# Initialize SparkSession
spark = SparkSession.builder \
    .appName("DataFrame SQL Example") \
    .getOrCreate()

# Sample data and schema
data = [("John", 28), ("Jane", 30), ("Jake", 25)]
columns = ["Name", "Age"]

# Create DataFrame
df = spark.createDataFrame(data, columns)

# Register DataFrame as a SQL temporary view
df.createOrReplaceTempView("people")

# Run SQL query
result = spark.sql("SELECT Name, Age FROM people WHERE Age > 25")

# Show the result
result.show()

# Stop the SparkSession
spark.stop()
```

ตัวอย่างการใช้ต่าง ๆ

```
from pyspark.sql import SparkSession

# Initialize a Spark session
spark = SparkSession.builder.master("local").appName("ReactionAnalysis").getOrCreate()

# Load the CSV files
file_path_1 = 'fb_live_thailand3.csv'
file_path_2 = 'fb_live_thailand2.csv'

# Read the data with header inference
df1 = spark.read.option("header", "true").csv(file_path_1)
df2 = spark.read.option("header", "true").csv(file_path_2)

combined_df = df1.union(df2)
combined_df.createOrReplaceTempView("reactions")
```

DISTINCT เป็นคำสั่งใน SQL (รวมถึง PySpark SQL) ที่ใช้ในการกรองข้อมูลเพื่อให้ได้เฉพาะค่าที่ไม่ซ้ำกันในคอลัมน์หรือชุดของคอลัมน์ที่ระบุ ซึ่งมีประโยชน์มากเมื่อคุณต้องการทราบค่าที่แตกต่างกันที่ปรากฏอยู่ในคอลัมน์ใดคอลัมน์หนึ่ง โดยไม่ต้องการให้ค่าซ้ำกัน

[+ Code](#) [+ Markdown](#)

```
result = spark.sql("SELECT DISTINCT num_reactions FROM reactions")
result.show()
```

Python

```
+-----+
|num_reactions|
+-----+
|          203|
|          221|
|          503|
|          227|
|          111|
|          135|
|          213|
|          150|
|          170|
|          217|
|          209|
|          313|
```

```
# Filter for num_reactions where the value is '262'
result = spark.sql("SELECT * FROM reactions WHERE status_type = 'video'")
result.show()
```

```
+-----+-----+-----+-----+-----+-----+
|Table|status_id|status_type|status_published|num_reactions|num_shares|
+-----+-----+-----+-----+-----+-----+
|FB3|246675545449582_1...|video|4/22/2018 6:00|529|262|
|FB3|246675545449582_1...|video|4/21/2018 6:17|227|57|
|FB3|246675545449582_1...|video|4/18/2018 0:24|503|72|
|FB3|246675545449582_1...|video|4/17/2018 7:42|295|53|
|FB3|246675545449582_1...|video|4/1/2018 5:16|332|30|
|FB3|246675545449582_1...|video|3/30/2018 8:28|135|79|
|FB3|246675545449582_1...|video|3/26/2018 8:28|150|47|
|FB3|246675545449582_1...|video|3/23/2018 7:09|221|36|
|FB2|246675545449582_1...|video|4/22/2018 6:00|529|512|
|FB2|246675545449582_1...|video|4/21/2018 6:17|227|236|
|FB2|246675545449582_1...|video|4/18/2018 0:24|503|614|
|FB2|246675545449582_1...|video|4/17/2018 7:42|295|453|
|FB2|246675545449582_1...|video|3/20/2018 1:28|18|0|
|FB2|246675545449582_1...|video|3/19/2018 22:34|37|0|
|FB2|246675545449582_1...|video|3/17/2018 8:07|75|36|
|FB2|246675545449582_1...|video|3/17/2018 7:47|90|78|
+-----+-----+-----+-----+-----+-----+
```

GROUP BY เป็นคำสั่งที่ใช้ใน SQL (รวมถึง PySpark SQL) เพื่อจัดกลุ่มข้อมูลตามค่าที่เหมือนกันในคอลัมน์ที่กำหนด ซึ่งจะช่วยให้เราสามารถสรุปข้อมูลได้สะดวกมากขึ้น เช่น การนับจำนวน การหาค่าเฉลี่ย หรือหาค่าสูงสุดต่ำสุดในแต่ละกลุ่ม

```
result = spark.sql("SELECT num_reactions, COUNT(*) as count FROM reactions GROUP BY num_reactions")
result.show()
```

Python

```
+-----+-----+
|num_reactions|count|
+-----+-----+
|          203|    2|
|          221|    1|
|          503|    2|
|          227|    3|
|          111|    2|
|          135|    1|
|          213|    2|
|          150|    3|
|          170|    1|
|          217|    2|
```

****ORDER BY ****ใช้ในการจัดเรียงข้อมูลตามลำดับที่ต้องการ เช่นเรียงจากน้อยไปมากหรือมากไปน้อย ASC: เรียงจากน้อยไปมาก (ค่าเริ่มต้น) DESC: เรียงจากมากไปน้อย

```
result = spark.sql("SELECT * FROM reactions ORDER BY num_reactions DESC")
result.show()
```

Table	status_id	status_type	status_published	num_reactions	num_shares
FB2	246675545449582_1...	photo	3/20/2018 1:54	98	0
FB2	246675545449582_1...	video	3/17/2018 7:47	90	78
FB2	246675545449582_1...	video	3/17/2018 8:07	75	36
FB3	246675545449582_1...	video	4/22/2018 6:00	529	262
FB2	246675545449582_1...	video	4/22/2018 6:00	529	512
FB3	246675545449582_1...	video	4/18/2018 0:24	503	72
FB2	246675545449582_1...	video	4/18/2018 0:24	503	614
FB2	246675545449582_1...	video	3/19/2018 22:34	37	0
FB3	246675545449582_1...	photo	4/5/2018 9:23	346	0

HAVING ใช้กับ GROUP BY เพื่กรองผลลัพธ์หลังจากทำการสรุปข้อมูลแล้ว โดยมีลักษณะการใช้งานคล้ายกับ WHERE แต่ใช้กับกลุ่มข้อมูลที่สรุปแล้ว

```
result = spark.sql("SELECT num_reactions, COUNT(*) as count FROM reactions GROUP BY num_reactions HAVING count > 1")
result.show()
```

Python

num_reactions	count
203	2
503	2
227	3
111	2
213	2
150	3
217	2
295	2
529	2

Filter PySpark

```
from pyspark.sql import SparkSession
```

```
spark = SparkSession.builder.appName("FilterExample").getOrCreate()
```

```
data = [(1, "Alice", 3000),
        (2, "Bob", 4500),
        (3, "Charlie", 5000),
        (4, "David", 3500)]
```

```
df = spark.createDataFrame(data, ["id", "name", "salary"])
```

```
# กรองพนักงานที่มีเงินเดือนมากกว่า 4000
```

```
filtered_df = df.filter(df.salary > 4000)
```

```
filtered_df.show()
```

```
# กรองข้อมูลที่มีเงินเดือนมากกว่า 3000 และน้อยกว่า 5000
filtered_df = df.filter((df.salary > 3000) & (df.salary < 5000))
filtered_df.show()

# กรองข้อมูลที่มีชื่อเป็น Bob หรือเงินเดือนมากกว่า 4000
filtered_df = df.filter((df.name == "Bob") | (df.salary > 4000))
filtered_df.show()
```

```
# กรองข้อมูลที่มีชื่อขึ้นต้นด้วย 'A'
filtered_df = df.filter(df.name.startswith("A"))
filtered_df.show()

# กรองข้อมูลที่มีชื่อมีตัวอักษร 'ar'
filtered_df = df.filter(df.name.like("%ar%"))
filtered_df.show()

# กรองข้อมูลที่มีชื่อมีลักษณะเป็น regex pattern
filtered_df = df.filter(df.name.rlike("a.*e")) # กรองชื่อที่ขึ้นต้นด้วย 'a' และลงท้ายด้วย 'e'
filtered_df.show()
```

```
# กรองพนักงานที่ id เป็น 1 หรือ 3
filtered_df = df.filter(df.id.isin([1, 3]))
filtered_df.show()
```

```
# กรองข้อมูลที่คอลัมน์ salary มีค่าไม่เป็น null
filtered_df = df.filter(df.salary.isNotNull())
filtered_df.show()

# กรองข้อมูลที่คอลัมน์ salary เป็น null
filtered_df = df.filter(df.salary.isNull())
filtered_df.show()
```

จาก Code Text Analytics

```
26 # กรองออกข้อมูลที่ ReviewText หรือ Rating เป็น null หรือว่างเปล่า
27 data = data.filter(
28     (col("ReviewText").isNotNull()) & (col("ReviewText") != "") &
29     (col("Rating").isNotNull()) & (col("Rating") != float('nan'))
30 )
```

```
# กรองข้อมูลใน DataFrame โดยใช้เงื่อนไขหลายประการ
data = data.filter(
    # เช็คค่าคอลัมน์ "ReviewText" มีค่า (ไม่เป็น null)
    (col("ReviewText").isNotNull()) &
    # เช็คค่าคอลัมน์ "ReviewText" ไม่ใช่ข้อความว่างเปล่า
    (col("ReviewText") != "") &
    # เช็คค่าคอลัมน์ "Rating" มีค่า (ไม่เป็น null)
    (col("Rating").isNotNull()) &
    # เช็คค่าคอลัมน์ "Rating" ไม่ใช่ค่า NaN (Not a Number)
    (col("Rating") != float('nan'))
)
```

ถ้าหากค่ามันว่างใน dataset

```
# แทนค่าว่างใน ReviewText ด้วย "No Review" และ Rating ด้วย 0
data = data.fillna({"ReviewText": "No Review", "Rating": 0.0})

# กรองข้อมูลที่เหลือ โดยไม่มี ReviewText เป็นข้อความว่าง
data = data.filter(
    (col("ReviewText") != "") &
    (col("Rating") != float('nan'))
)
```

```
from pyspark.sql.functions import col, mean

# คำนวณค่าเฉลี่ยของ Rating
avg_rating = data.select(mean(col("Rating"))).collect()[0][0]

# แทนที่ค่าว่างใน Rating ด้วยค่าเฉลี่ยที่คำนวณได้
data = data.fillna({"Rating": avg_rating})

# แทนค่าว่างใน ReviewText ด้วย "No Review"
data = data.fillna({"ReviewText": "No Review"})

data.show()
```

```
# ลบแถวที่มีค่า null หรือ NaN ในคอลัมน์ ReviewText หรือ Rating
data = data.dropna(subset=["ReviewText", "Rating"])

# ลบแถวที่ ReviewText เป็นข้อความว่างเปล่า
data = data.filter(col("ReviewText") != "")

data.show()
```