### DADS6005 MIDTERM REALTIME

Data streaming and Realtime analytics System Requirement a. Data source 3 sources

- a. Data source 3 sources [1]
- i. Source 1 : PageView (stream datagen) (topic1)

PageView เป็นแหล่งข้อมูลที่สร้างข้อมูลจำลองเกี่ยวกับการดูหน้าเว็บของผู้ใช้งาน โดยใช้ Kafka Connect Data Generator เพื่อสร้างข้อมูลในรูปแบบ JSON (ไม่มีการกำหนด schema) และส่งไปยัง Kafka topic 1\_pageviews ข้อมูลนี้ ประกอบด้วยเวลาในการดู (viewtime), รหัสผู้ใช้งาน (userid), และรหัสหน้าเว็บ (pageid) ซึ่งถูกสร้างขึ้นแบบสุ่มใน ช่วงเวลาที่กำหนด (random interval) จากนั้นใช้ Kafka Stream ชื่อ pageviews\_stream เพื่อดึงข้อมูลเข้าสู่ระบบ สำหรับการประมวลผลแบบเรียลไทม์ ข้อมูลนี้มีความสำคัญต่อการวิเคราะห์พฤติกรรมผู้ใช้ เช่น การติดตามการใช้งานเว็บไซต์ และการปรับปรุงกลยุทธ์การตลาด ตัวอย่างข้อมูลที่ได้ ได้แก่ {"viewtime": 691001, "userid": "User\_6", "pageid": "Page 22"} ซึ่งสามารถนำไปใช้งานในการสร้างรายงานหรือวิเคราะห์เชิงลึกได้

```
CREATE STREAM
                                                                                                "viewtime": 691001,
 "name": "datagen-pageviews",
                                                          pageviews_stream (
 "connector.class":
                                                                                                "userid": "User 6",
                                                           viewtime BIGINT,
                                                                                                "pageid": "Page_22"
"io.confluent.kafka.connect.datagen.DatagenConnector",
                                                           userid VARCHAR,
 "key.converter":
                                                           pageid VARCHAR
"org.apache.kafka.connect.storage.StringConverter",
                                                          )
 "value.converter":
                                                          WITH
"org.apache.kafka.connect.json.JsonConverter",
                                                          (KAFKA_TOPIC='1_pageviews',
                                                          VALUE_FORMAT='JSON');
 "value.converter.schemas.enable": "false",
 "kafka.topic": "1 pageviews",
 "max.interval": "1000",
 "quickstart": "pageviews",
 "interval.type": "random",
 "interval.range.min": "1",
 "interval.range.max": "100
```

## ii. Source 2 : Users\_ (stream datagen) (topic2)

Users เป็นแหล่งข้อมูลจำลองเกี่ยวกับผู้ใช้งานที่สร้างขึ้นโดย Kafka Connect Data Generator ข้อมูลนี้ส่งไปยัง Kafka topic 2\_users ในรูปแบบ JSON (ไม่มีการกำหนด schema) ประกอบด้วยเวลาลงทะเบียน (registertime), รหัส ผู้ใช้งาน (userid), ภูมิภาค (regionid), และเพศ (gender) โดยมีการสร้างข้อมูลในช่วงเวลาสุ่ม (random interval) จากนั้นสร้าง Kafka Table ชื่อ users\_tb ซึ่งกำหนด userid เป็น Primary Key เพื่อจัดเก็บข้อมูลดังกล่าว ข้อมูลนี้มี ความสำคัญต่อการวิเคราะห์พฤติกรรมผู้ใช้งานและเชื่อมโยงข้อมูลในระบบ เช่น การติดตามกลุ่มผู้ใช้งานในแต่ละภูมิภาค ตัวอย่างข้อมูล ได้แก่ {"registertime": 1491293306724, "userid": "User\_7", "regionid": "Region\_8", "gender": "MALE"} ซึ่งสามารถนำไปใช้ในการวิเคราะห์ข้อมูลเชิงลึกและการตลาดแบบเฉพาะกลุ่ม.

```
CREATE TABLE users_tb (
 "name": "datagen-users",
                                                           registertime BIGINT,
                                                                                               "registertime":
                                                           userid VARCHAR PRIMARY KEY,
"connector.class":
                                                                                             1491293306724,
"io.confluent.kafka.connect.datagen.DatagenConnector",
                                                           regionid VARCHAR,
                                                                                               "userid": "User_7",
"key.converter":
                                                           gender VARCHAR
                                                                                               "regionid": "Region_8",
"org.apache.kafka.connect.storage.StringConverter",
"value.converter":
                                                         WITH (KAFKA_TOPIC='2_users',
                                                                                               "gender": "MALE"
"org.apache.kafka.connect.json.JsonConverter",
                                                          VALUE_FORMAT='JSON');
"value.converter.schemas.enable": "false",
"kafka.topic": "2_users",
"max.interval": "1000",
"quickstart": "users",
"interval.type": "random",
"interval.range.min": "1",
"interval.range.max": "100"
```

#### .iii. Source 3 : Your design (relational database) (topic3)

Orders เป็นแหล่งข้อมูลที่จัดเก็บคำสั่งชื้อของผู้ใช้งานในรูปแบบ relational database โดยใช้ Python script สร้าง ข้อมูลจำลอง เช่น รหัสคำสั่งชื้อ (orderid), รหัสผู้ใช้งาน (userid), ประเภทสินค้า (product\_type), จำนวน (quantity), และสถานะ (status) จากนั้นส่งข้อมูลในรูปแบบ JSON ไปยัง Kafka topic 3\_orders ข้อมูลเหล่านี้ถูกนำเข้า Kafka Stream ชื่อ orders\_stream เพื่อประมวลผลแบบเรียลไทม์ ข้อมูลคำสั่งชื้อมีความสำคัญต่อการติดตามการขายสินค้า การ ประเมินยอดขาย และการปรับปรุงกระบวนการจัดส่ง ตัวอย่างข้อมูลที่สร้างขึ้น ได้แก่ {"ORDERID": "4e383b6b-3f38-45ad-baec-c100b17a4323", "USERID": "User\_7", "PRODUCT\_TYPE": "Tablet", "UNIT\_PRICE": 318.12, "QUANTITY": 3, "TOTAL\_PRICE": 954.36, "STATE": "Nebraska", "STATUS": "Shipped"} ซึ่งช่วยสนับสนุนการ วิเคราะห์ยอดขายและการจัดการคำสั่งชื้อในระบบ.

```
CREATE STREAM orders_stream (
  ORDERID STRING,
                                                           "ORDERID": "4e383b6b-3f38-45ad-baec-c100b17a4323",
  USERID STRING,
                                                           "USERID": "User 7",
  ORDER_TIMESTAMP STRING,
                                                           "ORDER TIMESTAMP": "2024-11-16T08:33:16.098003",
  PRODUCT_TYPE STRING,
                                                           "PRODUCT TYPE": "Tablet",
  UNIT_PRICE DOUBLE,
                                                           "UNIT PRICE": 318.12,
  QUANTITY INT,
                                                           "QUANTITY": 3,
  TOTAL PRICE DOUBLE,
                                                           "TOTAL PRICE": 954.36,
  STATE STRING,
                                                           "STATE": "Nebraska",
  STATUS STRING
                                                           "STATUS": "Shipped"
) WITH (
  KAFKA_TOPIC='3_orders',
  VALUE_FORMAT='JSON'
```

#### b. Kafka system [2]

#### i. 5 partitions

Kafka ใช้ partitions ในการแบ่งข้อมูลภายในแต่ละ topic เพื่อเพิ่มประสิทธิภาพในการจัดการและกระจายข้อมูล ไปยัง brokers หลายตัว ในระบบนี้ topic แต่ละอันถูกแบ่งออกเป็น 5 partitions ซึ่งช่วยให้รองรับการประมวลผลข้อมูล แบบขนาน (parallel processing) ได้ดีขึ้นและเพิ่มประสิทธิภาพทั้งในด้านการเขียนและการอ่านข้อมูล นอกจากนี้ การใช้ partitions ยังช่วยเพิ่มความสามารถในการขยายระบบ (scalability) และรองรับการใช้งานที่มีข้อมูลจำนวนมาก โดยข้อมูล จะถูกกระจายอย่างสมดุลไปยัง brokers ต่างๆ เพื่อให้ระบบทำงานได้อย่างราบรื่น.

#### ii. 3 brokers

ในการสร้างระบบ Kafka เราจะสร้าง 3 brokers หมายถึงมี 3 servers ที่ทำหน้าที่ในการจัดการและจัดเก็บข้อมูล โดย brokers จะรับข้อมูลจาก producer และส่งข้อมูลไปยัง consumer การใช้หลาย brokers ช่วยเพิ่มความทนทาน (fault tolerance) ของระบบ เนื่องจากหาก broker ตัวใดตัวหนึ่งล้มเหลว ข้อมูลยังคงสามารถเข้าถึงได้จาก brokers อื่นๆ ที่มีข้อมูลสำรอง (replica) อยู่ ทำให้ระบบมีความเสถียรและสามารถดำเนินการต่อไปได้โดยไม่กระทบต่อการทำงานของ ระบบ ในที่นี้ ระบบใช้ broker:29092, broker1:29095, และ broker2:29098 เพื่อจัดการข้อมูลและรองรับการกระจาย ข้อมูลระหว่าง servers.

### iii. 8 topics

ในการสร้าง 8 topics หมายถึงมีการตั้งค่าไว้ทั้งหมด 8 topic ที่ใช้ในการแยกประเภทข้อมูล โดยแต่ละ topic จะรับ ข้อมูลจาก producers และส่งข้อมูลไปยัง consumers เพื่อการประมวลผลหรือการใช้งานต่อไป การใช้หลายๆ topic ช่วย ให้สามารถจัดการและแยกประเภทข้อมูลที่แตกต่างกันได้อย่างมีประสิทธิภาพ เช่น หนึ่ง topic อาจใช้สำหรับการบันทึกข้อมูล การสั่งซื้อ (orders), อีก topic อาจใช้สำหรับการบันทึกข้อมูลผู้ใช้ (users), ซึ่งทำให้สามารถควบคุมข้อมูลแต่ละประเภทได้ อย่างเป็นระเบียบและสามารถเข้าถึงได้ง่ายขึ้นในระบบ

#### iv. Schema Register

Schema Registry เป็นเครื่องมือที่ช่วยในการจัดการ schema ของข้อมูลที่ถูกส่งไปยัง Kafka topics เช่น Avro, JSON, หรือ Protobuf โดยทำหน้าที่ตรวจสอบและรับรองว่า ข้อมูลที่ถูกส่งไปยัง Kafka มีรูปแบบที่ถูกต้องและสอดคล้องกับ schema ที่กำหนดไว้ ช่วยให้มั่นใจได้ว่าข้อมูลมีความสมบูรณ์และสามารถตรวจสอบได้เมื่อมีการส่งข้อมูลระหว่าง producers และ consumers ผ่าน Kafka โดยการใช้ Schema Registry ยังช่วยในการป้องกันปัญหาการเข้ากันไม่ได้ ระหว่างรูปแบบข้อมูลที่ต่างกัน เช่น เมื่อข้อมูลถูกส่งในรูปแบบ Avro หรือ JSON และช่วยให้ง่ายต่อการบำรุงรักษาและการ เปลี่ยนแปลง schema ในระบบได้อย่างราบรื่น.

#### v. Kafka connect

Kafka Connect เป็นเครื่องมือที่ช่วยเชื่อมต่อ Kafka กับแหล่งข้อมูลภายนอก เช่น ฐานข้อมูล, ระบบการจัดเก็บ ข้อมูล, หรือแอปพลิเคชันอื่นๆ โดยใช้ connectors ที่ช่วยในการดึงข้อมูลจากแหล่งภายนอกเข้าสู่ Kafka หรือส่งข้อมูลจาก Kafka ไปยังระบบภายนอก การใช้ Kafka Connect ช่วยให้การทำงานกับข้อมูลภายนอกเป็นไปอย่างสะดวกและรวดเร็ว โดยไม่จำเป็นต้องเขียนโค้ด เนื่องจาก connectors จะทำการเชื่อมต่อและถ่ายโอนข้อมูลอย่างอัตโนมัติ โดยในกรณีนี้เราใช้ Data Generator เพื่อสร้างข้อมูลจำลองเกี่ยวกับ pageviews และ users ซึ่งช่วยให้สามารถสร้างข้อมูลทดสอบได้ง่ายและ รวดเร็วโดยไม่ต้องใช้ข้อมูลจริง.

### c. ksqlDB operation [3]

### i. Clean or transform data (topic4)

เราได้ทำการ transforming ข้อมูลจากสองแหล่งข้อมูล (sources) คือ PAGEVIEWS\_STREAM และ users\_tb เราจะ ทำการแปลงข้อมูลที่มีอยู่ให้มีรูปแบบที่เข้าใจได้ง่ายและถูกต้องมากขึ้น:

- 1. PAGEVIEWS\_STREAM: ข้อมูลที่เกี่ยวข้องกับ viewtime หรือเวลาในการดูหน้าเว็บจะถูกแปลงจากหน่วย มิลลิวินาทีเป็นหน่วยนาที (โดยการหารด้วย 6000) เพื่อให้สามารถเข้าใจได้ง่ายขึ้น เช่น การแปลงข้อมูล VIEWTIME = 691001 มิลลิวินาที จะกลายเป็น VIEWTIME\_FORMATTED = 22 นาที ข้อมูลที่ถูกแปลงนี้จะถูก ส่งไปยัง Kafka topic 4\_cleaned โดยใช้คำสั่ง CREATE STREAM.
- 2. users\_tb: ข้อมูลที่เกี่ยวข้องกับ REGISTERTIME หรือเวลาลงทะเบียนผู้ใช้งานจะถูกแปลงจากมิลลิวินาทีเป็น รูปแบบเวลาที่เข้าใจได้ง่าย (เช่น yyyy-MM-dd HH:mm
- ) ด้วยการใช้ฟังก์ชัน TIMESTAMPTOSTRING ซึ่งจะทำให้เวลาที่มีรูปแบบเช่น 1970-01-18 17:30:02 สามารถอ่านและ ตีความได้ง่ายขึ้น ข้อมูลนี้จะถูกส่งไปยัง Kafka topic 4\_cleaned ด้วยคำสั่ง CREATE TABLE.

การแปลงข้อมูลเหล่านี้ช่วยให้ข้อมูลมีความถูกต้องและเข้าใจง่ายขึ้น ทำให้การประมวลผลข้อมูลในภายหลัง เช่น การ วิเคราะห์หรือการใช้งานในระบบอื่น ๆ เป็นไปได้อย่างมีประสิทธิภาพและมีความแม่นยำมากขึ้น.

```
CREATE STREAM PAGEVIEWS CLEAN WITH (KAFKA TOPIC =
'4_cleaned', VALUE_FORMAT = 'JSON') AS
                                                         "USERID": "User_2",
  SELECT
                                                         "PAGEID": "Page 97",
    USERID,
                                                         "VIEWTIME_FORMATTED": 22
    PAGEID,
        VIEWTIME/6000 AS VIEWTIME_FORMATTED
                                                       }
  FROM PAGEVIEWS_STREAM;
CREATE TABLE USERS CLEAN WITH (KAFKA TOPIC =
                                                       {
'4_cleaned', VALUE_FORMAT = 'JSON') AS
                                                         "USERID": "User_6",
  SELECT
                                                         "REGIONID": "Region 1",
    USERID,
                                                         "GENDER": "OTHER",
    REGIONID,
    GENDER,
                                                         "REGISTERTIME FORMATTED": "1970-01-18
    TIMESTAMPTOSTRING(REGISTERTIME / 1000, 'yyyy-MM-dd
                                                       17:30:02"
HH:mm:ss', 'Asia/Bangkok') AS REGISTERTIME_FORMATTED
                                                       }
  FROM users_tb;
```

### ii. Aggregation (join + group by) (topic5)

การ aggregation นี้ เราจะทำการรวบรวมข้อมูลคำสั่งชื้อจาก ORDERS\_STREAM และข้อมูลผู้ใช้จาก
USERS\_CLEAN โดยการใช้ JOIN ข้อมูลทั้งสองแหล่ง และทำการ GROUP BY ตามประเภทสินค้าที่สั่งชื้อ
(PRODUCT\_TYPE). การ JOIN จะทำให้เราสามารถเชื่อมโยงข้อมูลคำสั่งชื้อกับข้อมูลของผู้ใช้งานผ่าน USERID เพื่อให้ได้
ข้อมูลที่ครบถ้วนมากขึ้น.

หลังจากนั้น เราจะคำนวณข้อมูลต่าง ๆ ที่เกี่ยวข้อง เช่น:

- total\_orders: จำนวนคำสั่งซื้อทั้งหมดในแต่ละประเภทสินค้า
- total\_sales: ยอดขายรวมจากคำสั่งซื้อทั้งหมด
- avg\_sales\_per\_order: ยอดขายเฉลี่ยต่อคำสั่งซื้อ

ข้อมูลที่ได้จากการคำนวณนี้จะถูกส่งไปยัง Kafka topic 5\_aggregated ซึ่งจะช่วยให้ข้อมูลที่รวบรวมและประมวลผลแล้ว สามารถนำไปใช้งานหรือวิเคราะห์ได้ต่อไป เช่น การสร้างรายงาน หรือการใช้ข้อมูลในแอปพลิเคชันที่เกี่ยวข้อง

```
CREATE TABLE 5_sales_by_product
                                                        "PRODUCT TYPE": "Keyboard",
  WITH (KAFKA_TOPIC = '5_aggregated') AS
SELECT
                                                        "TOTAL ORDERS": 187,
  o.PRODUCT_TYPE,
                                                        "TOTAL_SALES": 49598.950000000026,
 COUNT(o.ORDERID) AS total_orders,
                                                        "AVG_SALES_PER_ORDER": 265.23502673796804
  SUM(o.TOTAL_PRICE) AS total_sales,
  SUM(o.TOTAL_PRICE) / COUNT(o.ORDERID) AS
                                                      }
avg_sales_per_order
FROM ORDERS STREAM o
LEFT JOIN USERS_CLEAN u
  ON o.USERID = u.USERID
GROUP BY o.PRODUCT_TYPE;
```

#### iii. Windows

### 1. Tumbling (topic6)

การ Tumbling Window นี้ เราจะใช้ฟังก์ชัน TUMBLING (SIZE 1 HOUR) เพื่อคำนวณจำนวนการดูหน้าเว็บ (pageviews) ของผู้ใช้แต่ละคนในช่วงเวลาที่เป็น 1 ชั่วโมง ซึ่งจะไม่ทับซ้อนกัน โดยการใช้ Tumbling Window จะ ช่วยแบ่งข้อมูลออกเป็นช่วงเวลาที่ชัดเจน เช่น ช่วงเวลา 1 ชั่วโมง ซึ่งข้อมูลแต่ละช่วงจะถูกคำนวณแยกจากกันและไม่ทับ ซ้อนกัน

การคำนวณจะทำการ COUNT จำนวนการดูหน้าเว็บ (pageviews) ของผู้ใช้ (USERID) ในแต่ละช่วงเวลา (window) และผลลัพธ์ที่ได้จะประกอบด้วย:

- USERID: รหัสผู้ใช้ที่ทำการดูหน้าเว็บ
- PAGEVIEWS COUNT: จำนวนครั้งที่ผู้ใช้งานดูหน้าเว็บในช่วงเวลา 1 ชั่วโมง
- WINDOWSTART และ WINDOWEND: เวลาเริ่มต้นและสิ้นสุดของช่วงเวลา 1 ชั่วโมงที่ใช้ในการคำนวณ ข้อมูลที่ได้จะถูกส่งไปยัง Kafka topic 6\_tumbling ซึ่งสามารถนำไปใช้งานต่อในระบบการประมวลผลข้อมูลแบบ เรียลไทม์ หรือใช้ในการวิเคราะห์พฤติกรรมของผู้ใช้ในช่วงเวลาต่าง ๆ.

```
CREATE TABLE 6_pageviews_per_user WITH

(KAFKA_TOPIC='6_tumbling') AS

SELECT

pv.USERID,

COUNT(*) AS PAGEVIEWS_COUNT

FROM PAGEVIEWS_CLEAN pv

WINDOW TUMBLING (SIZE 1 HOUR) -- 1-hour window

GROUP BY pv.USERID EMIT CHANGES;

{

"USERID": "User_7",

"WINDOWSTART": 1731740400000,

"WINDOWSTART": 1731744000000,

"WINDOWEND": 1731744000000,

"PAGEVIEWS_COUNT": 153

}
```

## 2. Hopping (topic7)

การ Hopping Window นี้ เราจะใช้ HOPPING (SIZE 1 HOUR, ADVANCE BY 30 MINUTES) เพื่อคำนวณ จำนวนคำสั่งชื้อและยอดขายรวมของผู้ใช้แต่ละคนในช่วงเวลา 1 ชั่วโมง โดยการเลื่อนช่วงเวลา (hop) ทุก ๆ 30 นาที. การใช้ Hopping Window ช่วยให้สามารถคำนวณข้อมูลได้ในหลาย ๆ ช่วงเวลา โดยที่ทุก ๆ 30 นาทีจะมีการเลื่อน ช่วงเวลาใหม่เข้ามา ซึ่งทำให้สามารถคำนวณข้อมูลใหม่ในแต่ละช่วงเวลาได้อย่างต่อเนื่องและซ้ำข้อนกัน. ผลลัพธ์ที่ได้จากการคำนวณจะประกอบด้วย:

- USERID: รหัสผู้ใช้ที่ทำการสั่งซื้อ
- ORDER\_COUNT: จำนวนคำสั่งซื้อของผู้ใช้ในช่วงเวลานั้น
- TOTAL\_SALES: ยอดขายรวมจากคำสั่งซื้อในช่วงเวลานั้น
- WINDOWSTART และ WINDOWEND: เวลาเริ่มต้นและสิ้นสุดของช่วงเวลา (window) ที่ใช้ในการคำนวณ ข้อมูลที่ได้จะถูกส่งไปยัง Kafka topic 7\_hopping ซึ่งสามารถใช้ในการวิเคราะห์และติดตามพฤติกรรมการสั่งซื้อของ ผู้ใช้ในช่วงเวลาต่าง ๆ เช่น การวิเคราะห์ยอดขายในช่วงเวลาที่มีการเปลี่ยนแปลงอย่างรวดเร็ว.

```
CREATE TABLE 7_total_orders_hopping

WITH (KAFKA_TOPIC='7_hopping') AS

SELECT

o.USERID,

COUNT(*) AS ORDER_COUNT,

SUM(o.TOTAL_PRICE) AS TOTAL_SALES

FROM orders_stream o

WINDOW HOPPING (SIZE 1 HOUR, ADVANCE BY 30 MINUTES)

-- 1-hour window, hops every 30 minutes

GROUP BY o.USERID EMIT CHANGES;

{

"USERID": "User_7",

"WINDOWSTART": 1731740400000,

"WINDOWSTART": 1731744000000,

"WINDOWEND": 1731744000000,

"PAGEVIEWS_COUNT": 153

}
```

#### 3. Session (topic8)

การ Session Window นี้ เราจะใช้ Session Window เพื่อติดตามและคำนวณจำนวนคำสั่งชื้อและยอดรวมคำสั่ง ชื้อของผู้ใช้แต่ละคนในช่วงเวลาที่กำหนด 10 นาที ซึ่งระบบจะพิจารณาว่าเซสชันสิ้นสุดเมื่อไม่มีการกระทำจากผู้ใช้ ภายในช่วงเวลา 10 นาที. เมื่อไม่มีคำสั่งซื้อใหม่เกิดขึ้นในช่วงเวลานั้น เซสชันจะถูกปิดและเริ่มเซสชันใหม่ทันที. ผลลัพธ์ที่ได้จะประกอบด้วย:

- ORDER\_COUNT: จำนวนคำสั่งซื้อที่ผู้ใช้ทำในเซสซันนั้น
- TOTAL\_ORDER\_VALUE: ยอดรวมของคำสั่งซื้อที่เกิดขึ้นในเซสชัน
- WINDOWSTART และ WINDOWEND: เวลาที่เริ่มต้นและสิ้นสุดของเซสชันนั้น

ข้อมูลที่ได้จะถูกส่งไปยัง Kafka topic 8\_session, ซึ่งสามารถนำไปใช้ในการวิเคราะห์พฤติกรรมการสั่งซื้อของผู้ใช้ที่มี การกระทำหรือกิจกรรมอย่างต่อเนื่องในช่วงเวลาที่กำหนด และช่วยในการติดตามความสัมพันธ์ระหว่างผู้ใช้และการ สั่งซื้อในช่วงเวลานั้น ๆ.

```
CREATE TABLE 8_orders_session_table
                                                     {
  WITH (KAFKA_TOPIC='8_session') AS
                                                      "USERID": "User 5",
  SELECT
                                                      "WINDOWSTART": 1731745993095,
    USERID.
                                                      "WINDOWEND": 1731747620802,
    COUNT(*) AS ORDER COUNT,
    SUM(TOTAL_PRICE) AS TOTAL_ORDER_VALUE
                                                      "ORDER_COUNT": 180,
  FROM orders_stream
                                                      "TOTAL_ORDER_VALUE": 278670.22
  WINDOW SESSION (10 MINUTES)
                                                     }
  GROUP BY USERID
  EMIT CHANGES;
```

### iv. Testing the correctness [4]

สำหรับการทำ testing correctness ในระบบที่ใช้ ksqlDB, การทดสอบจะมุ่งเน้นที่การตรวจสอบความถูกต้องของ การประมวลผลข้อมูลและการทำงานของ streams และ tables ที่ถูกสร้างขึ้นใน ksqlDB เพื่อให้แน่ใจว่าข้อมูลที่ได้จากการ คิวรีเป็นไปตามที่คาดหวังและสามารถทำงานได้ตามความต้องการ. โดยในที่นี้เราจะทำการสร้าง input.json, output.json และ statement.sql ในการทดสอบ

การทดสอบความถูกต้องใน ksqlDB มุ่งเน้นที่การตรวจสอบการประมวลผลข้อมูลจาก Kafka stream ว่าถูกต้อง หรือไม่ โดยใช้ input.json, output.json, และ statement.sql ดังนี้:

- 1. input.json: ข้อมูลตัวอย่างที่ส่งไปยัง Kafka topic เช่น USER\_ID, PAGE\_URL, และ TIMESTAMP.
- 2. statement.sql: คำสั่ง SQL สำหรับสร้าง stream และคิวรีข้อมูลจาก Kafka topic.
- 3. output.json: ผลลัพธ์ที่คาดหวังจากการคิวรีข้อมูลใน ksqlDB.

โดยข้อมูลใน input.json จะถูกส่งไป Kafka, จากนั้นจะใช้คำสั่งใน statement.sql เพื่อดึงข้อมูลและเปรียบเทียบกับ output.json เพื่อทดสอบความถูกต้องในการประมวลผลข้อมูล.

```
CREATE STREAM pageviews (
USER_ID STRING,
PAGE_URL STRING,
TIMESTAMP TIMESTAMP
) WITH (
KAFKA_TOPIC='datagen-pageviews',
VALUE_FORMAT='JSON'
);

SELECT USER_ID, PAGE_URL, TIMESTAMP
FROM pageviews
LIMIT 10; — LIMIT to check a small
batch of records
```

```
Get: Whttps://download.docker.com/linus/ubuntu_jammy_InRelease [48.8 kB]
Get: Shttp://ap-southeast-l.ec2.archive.ubuntu.com/ubuntu_jammy_updates/main_amd6W Packages [2151 kB]
Get: Shttp://ap-southeast-l.ec2.archive.ubuntu.com/ubuntu_jammy_updates/universe_amd6W Packages [1135 kB]
Get: Shttp://ap-southeast-l.ec2.archive.ubuntu.com/ubuntu_jammy_updates/universe_amd6W Packages [1135 kB]
Get: Thttp://security.ubuntu.com/ubuntu_jammy_security InRelease [129 kB]
Fetched 7718 kB in 1s [2583 kB/s]
Reading package lists...Done
ubuntuBja-172-31-4W-86:-5 sudo apt-get install confluent-cli
Reading package lists...Done
Reading state information...Done
Reading state information...Done
Reading state information...Done
Command 'taq': from deb freetds-bin (1.3.6-1)
command 'taq': from deb freetds-bin (1.3.6-1)
command 'saq': from deb freetds-bin (1.3.6-1)
command 'saq': from deb parallel (2021802345-2)
command 'saq': from deb parallel (2021802345-2)
command 'saq': from deb parallel (2021802345-2)
Try: sudo apt install 'cdeb name'
ubuntuBja-172-31-4W-86:-$ sudo nano statements.sql
ubuntuBja-172-31-4W-86:-$ docker or phome/ubuntu/statements.sql
ubuntuBja-172-31-4W-86:-$ docker or phome/ubuntu/statements.sql
ubuntuBja-172-31-4W-86:-$ docker or phome/ubuntu/statements.sql
ubuntuBja-172-31-4W-86:-$ docker or care the Hebb7ad2965 /bin/base/apuser/statements.sql
ubuntuBja-172-31-4W-86:-$ docker or care the Hebb7ad2965 /bin/base/apuser/statements.sql
ibuntuBja-172-31-4W-86:-$ docker or care the Hebb7ad2965 /bin/
```

# d. Apache Pinot [5,6]

Apache Pinot เป็นระบบจัดเก็บข้อมูลแบบคิวรีเรียลไทม์ที่ออกแบบมาเพื่อให้สามารถประมวลผลข้อมูลจาก แหล่งข้อมูลภายนอก เช่น Kafka และทำการคิวรีข้อมูลแบบเรียลไทม์ได้อย่างมีประสิทธิภาพ เราได้ใช้ Apache Pinot ในการ กำหนด schema และสร้าง real-time tables เพื่อจัดการกับข้อมูลที่ถูกนำเข้ามาจาก Kafka ซึ่งจะช่วยให้เราสามารถ ประมวลผลข้อมูลในเวลาจริงและทำการคิวรีข้อมูลได้ทันที เราได้สร้าง 3 ตารางใน Apache Pinot ซึ่งมาจาก 3 แหล่งข้อมูล หลัก ได้แก่ pageviews, users, และ orders การตั้งค่าระบบนี้ช่วยให้เราสามารถดึงข้อมูลจาก Kafka มาประมวลผลใน ระบบของ Apache Pinot และใช้ข้อมูลเหล่านั้นเพื่อการวิเคราะห์ที่รวดเร็วและมีประสิทธิภาพในเวลาจริง.

#### [i]. At least 3 queries

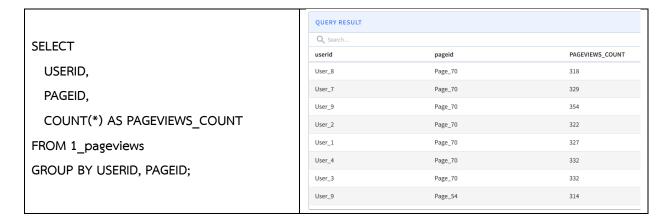
- คิวรีข้อมูลเฉลี่ยราคาต่อรัฐ : ในการคิวรีข้อมูลเฉลี่ยราคาต่อรัฐ เราจะคำนวณค่าเฉลี่ยของราคาสินค้า
   (AVG(TOTAL\_PRICE)) แยกตามรัฐ (STATE) จากข้อมูลที่จัดเก็บในตาราง 3\_orders ซึ่งบันทึกการสั่งชื้อทั้งหมด
   โดยจะใช้การจัดกลุ่มข้อมูลตามรัฐ (GROUP BY STATE) เพื่อคำนวณค่าเฉลี่ยราคาสินค้าของแต่ละรัฐ. ผลลัพธ์ที่ได้
   จะประกอบด้วย:
- STATE: รหัสหรือชื่อของรัฐที่มีการสั่งซื้อ
- avg\_price\_per\_state: ค่าเฉลี่ยของราคาสินค้าต่อรัฐค่าเฉลี่ยของราคาสินค้าต่อรัฐ

ข้อมูลที่ได้จะถูกใช้ในการวิเคราะห์และตัดสินใจทางธุรกิจ เช่น การกำหนดราคาและกลยุทธ์การตลาดที่เหมาะสมตามลักษณะ ของแต่ละรัฐ โดยผลลัพธ์จะถูกส่งไปยัง dashboard เพื่อแสดงผลต่อไป

**QUERY RESULT** Q Search.. STATE avg\_price\_per\_state **SELECT** Missouri 1529.3299882766712 STATE, Indiana 1545.0907540173052 AVG(TOTAL\_PRICE) AS Illinois 1651.6371727748694 avg\_price\_per\_state Michigan FROM 3\_orders Minnesota 1626.2678299776283 **GROUP BY STATE:** Kansas 1615.9134561626431 Ohio 1497.04423030303 1617.5715301204816 Iowa

- 2. คิวรีจำนวนการดูหน้าเว็บต่อผู้ใช้และหน้าเว็บ : คำนวณจำนวนการดูหน้าเว็บ (PAGEVIEWS\_COUNT) สำหรับแต่ ละผู้ใช้ (USERID) และแต่ละหน้าเว็บ (PAGEID) จากข้อมูลในตาราง 1\_pageviews ซึ่งบันทึกข้อมูลการดูหน้าเว็บ ทั้งหมด โดยใช้การจัดกลุ่มข้อมูลตาม USERID และ PAGEID. ผลลัพธ์ที่ได้จะประกอบด้วย:
- USERID: รหัสผู้ใช้ที่ดูหน้าเว็บ
- PAGEID: รหัสหรือชื่อของหน้าเว็บที่ถูกดู
- PAGEVIEWS\_COUNT: จำนวนการดูหน้าเว็บของผู้ใช้ในแต่ละหน้าเว็บ

ข้อมูลที่ได้จากคิวรีนี้จะถูกนำไปใช้ในการวิเคราะห์พฤติกรรมการเข้าชมของผู้ใช้ โดยสามารถดูได้ว่าผู้ใช้แต่ละคนดูหน้าเว็บใด บ่อยที่สุด หรือหน้าเว็บไหนที่ได้รับความนิยมมากที่สุด ซึ่งข้อมูลนี้สามารถช่วยในการปรับปรุงประสบการณ์ผู้ใช้บนเว็บไซต์ และสนับสนุนการตัดสินใจทางการตลาด.



- 3. คำนวณยอดขายรวม (SUM(TOTAL\_PRICE)) ตามประเภทของสินค้า (PRODUCT\_TYPE) จาก table 3\_orders ซึ่งบันทึกข้อมูลการสั่งซื้อ โดยจะแสดงผลลัพธ์ยอดขายรวมของแต่ละประเภทสินค้า ผลลัพธ์ที่ได้:
  - PRODUCT\_TYPE: ประเภทของสินค้าที่ขาย
  - TOTAL\_SALES: ยอดขายรวมของแต่ละประเภทสินค้า

ข้อมูลที่ได้จากคิวรีนี้จะช่วยในการวิเคราะห์ยอดขายตามประเภทสินค้า เช่น เพื่อตรวจสอบสินค้าที่มียอดขายสูงที่สุด หรือ สินค้าที่มีแนวโน้มในการขายดีในช่วงเวลาหนึ่ง ๆ ซึ่งสามารถนำไปใช้ในการวางแผนการจัดสรรสินค้าและการตัดสินใจทาง การตลาด.

SELECT

PRODUCT\_TYPE,

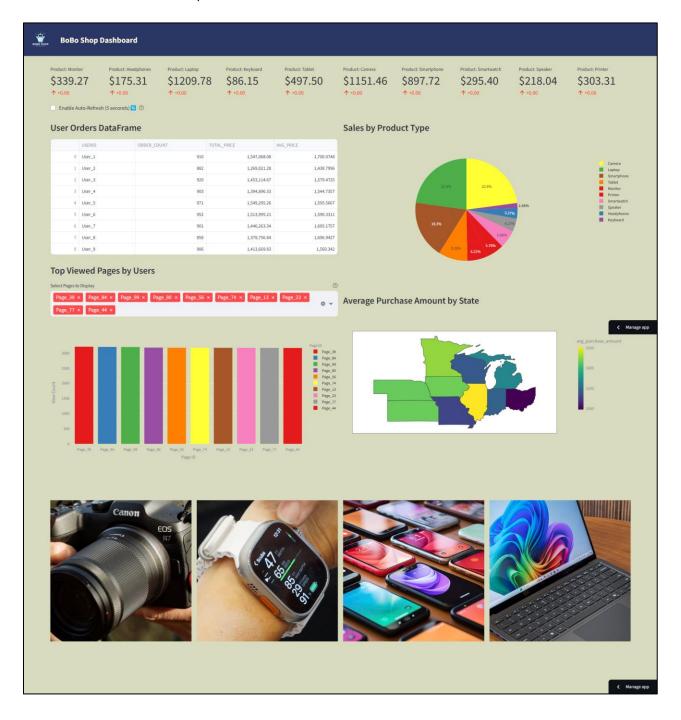
SUM(TOTAL\_PRICE) AS total\_sales

FROM 3\_orders

GROUP BY PRODUCT\_TYPE;

QUERY RESULT  Q. Search		
User_8	Page_70	318
User_7	Page_70	329
User_9	Page_70	354
User_2	Page_70	322
User_1	Page_70	327
User_4	Page_70	332
User_3	Page_70	332
User_9	Page_54	314

### e. Dashboard [7] i. At least 4 panels ii. Bonus: Interactive dashboard



ในส่วนของการแสดงผลข้อมูลบน Streamlit: <u>BoBoShop,</u> เราได้พัฒนา dashboard เพื่อให้สามารถดูข้อมูลต่างๆ ที่มีการอัปเดตในแบบ เรียลไทม์ โดยการเชื่อมต่อกับ Apache Pinot ผ่านพอร์ต 8099 เพื่อดึงข้อมูลจากฐานข้อมูลและแสดงผลได้ดังนี้:

# 1. ค่าเฉลี่ยราคาต่อประเภทสินค้า (Product Type)

ใช้ st.metric() เพื่อแสดงค่าเฉลี่ยราคาของสินค้าแต่ละประเภท โดยคำนวณจากคำสั่ง SQL ที่ดึงข้อมูลจากตาราง 3\_orders. ผลลัพธ์ จะถูกแสดงในรูปแบบของ metric สำหรับแต่ละประเภทสินค้า พร้อมทั้งแสดงความแตกต่างระหว่างข้อมูลล่าสุดกับข้อมูลก่อนหน้า เพื่อให้ผู้ใช้สามารถติดตามการเปลี่ยนแปลงของราคาสินค้าได้อย่างชัดเจนและทันที.

# 2. พฤติกรรมการสั่งซื้อของลูกค้า (Behavior User Data Frame)

ใช้ st.dataframe() เพื่อแสดงข้อมูลพฤติกรรมการสั่งซื้อของลูกค้า ซึ่งจะแสดงข้อมูลเช่น จำนวนการสั่งซื้อ, ยอดรวม, และ ค่าเฉลี่ย ของการสั่งซื้อ. ฟังก์ชันนี้ช่วยให้เราสามารถดูข้อมูลเชิงลึกในรูปแบบที่มีระเบียบและเข้าใจง่าย โดยการแสดงผลในรูปแบบ dataframe ที่อัปเดตแบบเรียลไทม์.

#### 3. ยอดขายตามประเภทสินค้า

ใช้ plotly.express.pie เพื่อแสดง ยอดขายทั้งหมด แยกตามประเภทสินค้าในรูปแบบ Pie Chart. กราฟนี้จะช่วยให้ผู้ใช้เห็นภาพรวม ของยอดขายในแต่ละประเภทสินค้าได้อย่างชัดเจนและเข้าใจง่าย โดยการใช้แผนภาพที่เข้าใจได้ง่ายและมีสีสันที่โดดเด่นเพื่อแสดง สัดส่วนยอดขายแต่ละประเภท.

### 4. จำนวนการดูหน้าเว็บตามผู้ใช้

ข้อมูลจำนวนการดูหน้าเว็บ (view\_count) จะถูกแสดงตาม PAGEID ของแต่ละหน้าเว็บที่ผู้ใช้เข้าชม โดยสามารถกรองข้อมูลตามหน้า เว็บที่เลือกได้จาก multiselect. ฟังก์ชันนี้ช่วยให้ผู้ใช้สามารถเลือกหน้าเว็บที่ต้องการดูข้อมูลและแสดงกราฟที่อัปเดตตามการเลือก, โดยการติดตามข้อมูลการเข้าชมแต่ละหน้าเว็บจะทำได้ง่ายและสะดวก.

# ยอดการซื้อเฉลี่ยตามรัฐ

ใช้ plotly.express.choropleth เพื่อแสดง ยอดการซื้อเฉลี่ย (avg\_purchase\_amount) ตามรัฐต่างๆ ของสหรัฐอเมริกาในรูปแบบ Choropleth map. แผนที่จะแบ่งสีตามค่าของยอดการซื้อที่แตกต่างกันในแต่ละรัฐ โดยการใช้แผนที่แบบสีที่มีการเน้นให้เห็นความ แตกต่างของยอดการซื้อระหว่างรัฐต่างๆ ช่วยให้ผู้ใช้เห็นข้อมูลการซื้อที่แตกต่างกันในแต่ละพื้นที่ได้อย่างซัดเจน.

#### ฟีเจอร์เสริม:

- Auto-refresh: หน้า dashboard มีการแสดงผลแบบอัปเดตอัตโนมัติทุก 5 วินาทีผ่าน checkbox เพื่อให้ข้อมูลที่แสดงผลมีความ ทันสมัยและถูกต้องตลอดเวลา.
- การแสดงภาพสินค้า: ด้านท้ายของ dashboard ยังมีการแสดง รูปภาพสินค้า เพื่อเสริมการนำเสนอข้อมูลและทำให้การนำเสนอข้อมูล
   มีความหลากหลายและน่าสนใจมากยิ่งขึ้น.

การใช้งาน Streamlit ร่วมกับ Apache Pinot ทำให้สามารถแสดงข้อมูลจาก Kafka แบบเรียลไทม่ได้อย่างมีประสิทธิภาพ โดยสามารถติดตาม การเปลี่ยนแปลงและแสดงผลข้อมูลเชิงลึกได้สะดวกผ่าน dashboard ที่ใช้งานง่าย.