고객을 세그먼테이션하자 [프로젝트] 정선아

11-2. 데이터 불러오기

데이터 살펴보기

• 테이블에 있는 10개의 행만 출력하기

```
SELECT *
FROM proven-serenity-439401-f2.modulabs_project.data
LIMIT 10
```



• 전체 데이터는 몇 행으로 구성되어 있는지 확인하기

```
SELECT COUNT(*) AS total_rows
FROM proven-serenity-439401-f2.modulabs_project.data
```



데이터 수 세기

• COUNT 함수를 사용해서, 각 컬럼별 데이터 포인트의 수를 세어 보기

```
SELECT

COUNT(InvoiceNo) AS COUNT_InvoiceNo,
COUNT(StockCode) AS COUNT_StockCode,
COUNT(Description) AS COUNT_Description,
COUNT(Quantity) AS COUNT_Quantity,
COUNT(InvoiceDate) AS COUNT_InvoiceDate,
COUNT(UnitPrice) AS COUNT_UnitPrice,
COUNT(CustomerID) AS COUNT_CustomerID,
COUNT(Country) AS COUNT_Country
FROM
proven-serenity-439401-f2.modulabs_project.data
```



11-4. 데이터 전처리 방법(1): 결측치 제거

컬럼 별 누락된 값의 비율 계산

- 각 컬럼 별 누락된 값의 비율을 계산
 - \circ 각 컬럼에 대해서 누락 값을 계산한 후, 계산된 누락 값을 UNION ALL을 통해 합치기

```
-- InvoiceNo에 대한 결측치 비율 계산
SELECT
    'InvoiceNo' AS column_name,
   ROUND(SUM(CASE WHEN InvoiceNo IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentage
   proven-serenity-439401-f2.modulabs_project.data
UNION ALL
-- StockCode에 대한 결측치 비율 계산
SELECT
    'StockCode' AS column_name,
   ROUND(SUM(CASE WHEN StockCode IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentage
FROM
   proven-serenity-439401-f2.modulabs_project.data
UNION ALL
-- Description에 대한 결측치 비율 계산
SELECT
    'Description' AS column_name,
   ROUND(SUM(CASE WHEN Description IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentag
FROM
   proven-serenity-439401-f2.modulabs_project.data
UNION ALL
-- Quantity에 대한 결측치 비율 계산
    'Quantity' AS column_name,
    ROUND(SUM(CASE WHEN Quantity IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentage
   proven-serenity-439401-f2.modulabs_project.data
UNION ALL
-- InvoiceDate에 대한 결측치 비율 계산
SELECT
    'InvoiceDate' AS column_name,
   ROUND(SUM(CASE WHEN InvoiceDate IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentag
   proven-serenity-439401-f2.modulabs_project.data
UNION ALL
-- UnitPrice에 대한 결측치 비율 계산
    'UnitPrice' AS column name,
   ROUND(SUM(CASE WHEN UnitPrice IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentage
   proven-serenity-439401-f2.modulabs_project.data
UNION ALL
-- CustomerID에 대한 결측치 비율 계산
SELECT
    'CustomerID' AS column_name,
   ROUND(SUM(CASE WHEN CustomerID IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentage
```

```
FROM proven-serenity-439401-f2.modulabs_project.data
UNION ALL
-- Country에 대한 결측치 비율 계산
SELECT
 'Country' AS column_name,
 ROUND(SUM(CASE WHEN Country IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100, 2) AS missing_percentage
FROM
 proven-serenity-439401-f2.modulabs_project.data;
```

| 작업 정 | 성보 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
|------|---------------|----|------|-------------|--------|
| 행 | column_name ▼ | // | | _percentage | |
| 1 | Description | | | 0.27 | |
| 2 | StockCode | | | 0.0 | |
| 3 | Country | | | 0.0 | |
| 4 | InvoiceDate | | | 0.0 | |
| 5 | UnitPrice | | | 0.0 | |
| 6 | CustomerID | | | 24.93 | |
| 7 | Quantity | | | 0.0 | |
| 8 | InvoiceNo | | | 0.0 | |

결측치 처리 전략

• StockCode = '85123A' 의 Description 을 추출하는 쿼리문을 작성하기

```
SELECT DISTINCT Description
FROM proven-serenity-439401-f2.modulabs_project.data
WHERE StockCode = '85123A';
```

쿼리 결과



결측치 처리

• DELETE 구문을 사용하며, WHERE 절을 통해 데이터를 제거할 조건을 제시

DELETE FROM proven-serenity-439401-f2.modulabs_project.data WHERE Description IS NULL OR CustomerID IS NULL;



11-5. 데이터 전처리(2): 중복값 처리

중복값 확인

- 중복된 행의 수를 세어보기
 - 。 8개의 컬럼에 그룹 함수를 적용한 후, COUNT가 1보다 큰 데이터를 세어보기

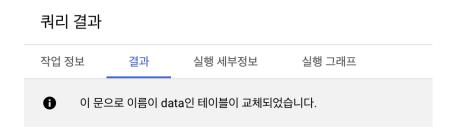
```
SELECT COUNT(*) AS duplicate_count
FROM (
 SELECT
   InvoiceNo,
   StockCode,
   Description,
   Quantity,
   InvoiceDate,
   UnitPrice,
   CustomerID,
   Country,
   COUNT(*) AS count_per_row
   proven-serenity-439401-f2.modulabs_project.data
   InvoiceNo, StockCode, Description, Quantity, InvoiceDate, UnitPrice, CustomerID, Country
 HAVING COUNT(*) > 1 -- 중복된 행 필터링
) AS duplicate_rows;
```

쿼리 결과 차트 JSON 실행 세부정보 실행 그래프 행 duplicate_count ▼ 1 4837

중복값 처리

- 중복값을 제거하는 쿼리문 작성하기
 - CREATE OR REPLACE TABLE 구문을 활용하여 모든 컬럼(*)을 DISTINCT 한 데이터로 업데이트

CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.data AS
SELECT DISTINCT *
FROM proven-serenity-439401-f2.modulabs_project.data ;



11-6. 데이터 전처리(3): 오류값 처리

InvoiceNo 살펴보기

• 고유(unique)한 InvoiceNo 의 개수를 출력하기

SELECT COUNT(DISTINCT InvoiceNo) AS unique_invoice_count FROM proven-serenity-439401-f2.modulabs_project.data



• 고유한 InvoiceNo 를 앞에서부터 100개를 출력하기

SELECT InvoiceNo AS unique_invoice FROM proven-serenity-439401-f2.modulabs_project.data LIMIT 100;

| 작업 정보 | 보 결과 | 차트 | JSON | 실행 세부정보 | 실행 그 |
|-------|----------------|----|------|---------|------|
| 1 | unique_invoice | ▼ | / | | |
| 88 | 549735 | | | | |
| 89 | 554032 | | | | |
| 90 | 561387 | | | | |
| 91 | 561387 | | | | |
| 92 | 561387 | | | | |
| 93 | 561387 | | | | |
| 94 | 561387 | | | | |
| 95 | 561387 | | | | |
| 96 | 561387 | | | | |
| 97 | 561387 | | | | |
| 98 | 561387 | | | | |
| 99 | 561387 | | | | |
| 100 | 574868 | | | | |

• InvoiceNo 가 'C'로 시작하는 행을 필터링 할 수 있는 쿼리문을 작성하기 (100행까지만 출력)

SELECT *
FROM proven-serenity-439401-f2.modulabs_project.data
WHERE InvoiceNo LIKE 'C%'
LIMIT 100;



• 구매 건 상태가 Canceled 인 데이터의 비율(%) - 소수점 첫번째 자리까지

SELECT ROUND(SUM(CASE WHEN InvoiceNo LIKE 'C%' THEN 1 ELSE 0 END)/COUNT(*) * 100, 1) AS canceled_percent FROM proven-serenity-439401-f2.modulabs_project.data

쿼리 결과



StockCode 살펴보기

• 고유한 StockCode 의 개수를 출력하기

SELECT COUNT(DISTINCT StockCode) AS unique_stockcode_count FROM proven-serenity-439401-f2.modulabs_project.data

쿼리 결과

| 작업 정 | 보 | 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
|------|-------|-------------|-----------|------|---------|--------|
| 행 // | uniqu | e_stockcode | e_count 🔻 | | | |
| 1 | | | 3684 | | | |

- 어떤 제품이 가장 많이 판매되었는지 보기 위하여 StockCode 별 등장 빈도를 출력하기
 - 。 상위 10개의 제품들을 출력하기

```
SELECT StockCode, COUNT(*) AS sell_cnt
FROM proven-serenity-439401-f2.modulabs_project.data
GROUP BY StockCode
ORDER BY sell_cnt DESC
LIMIT 10;
```

쿼리 결과

| 작업 정 | 보 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
|------|-------------|----|------------|---------|--------|
| 행 // | StockCode ▼ | // | sell_cnt ▼ | // | |
| 1 | 85123A | | | 2065 | |
| 2 | 22423 | | | 1894 | |
| 3 | 85099B | | | 1659 | |
| 4 | 47566 | | | 1409 | |
| 5 | 84879 | | | 1405 | |
| 6 | 20725 | | | 1346 | |
| 7 | 22720 | | | 1224 | |
| 8 | POST | | | 1196 | |
| 9 | 22197 | | | 1110 | |
| 10 | 23203 | | | 1108 | |

- StockCode 의 컬럼에 있던 값 중에서 숫자를 제외한 문자만 남기고 문자가 몇 자리 수 인지 세고
 - 。 **숫자가 0~1개인 값**들에는 어떤 코드들이 들어가 있는지 출력하기

```
SELECT DISTINCT StockCode, number_count

FROM (
SELECT StockCode,

LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS number_count
```

```
FROM proven-serenity-439401-f2.modulabs_project.data
) AS stock_with_number_count
WHERE number_count <= 1;
```

| 작업 정 | 보 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
|------|--------------|----|------------|---------|--------|
| 행 // | StockCode ▼ | // | number_cou | unt ▼ | |
| 1 | POST | | | 0 | |
| 2 | М | | | 0 | |
| 3 | PADS | | | 0 | |
| 4 | D | | | 0 | |
| 5 | BANK CHARGES | | | 0 | |
| 6 | DOT | | | 0 | |
| 7 | CRUK | | | 0 | |
| 8 | C2 | | | 1 | |

- StockCode 의 컬럼에 있던 값 중에서 숫자를 제외한 문자만 남기고 문자가 몇 자리 수 인지 세고
 - 숫자가 0~1개인 값들을 가지고 있는 데이터 수는 전체 데이터 수 대비 몇 퍼센트인지 구하기 (소수점 두 번째 자리까지)

```
WITH StockCodesWithFewNumbers AS (
SELECT StockCode
FROM `proven-serenity-439401-f2.modulabs_project.data`
WHERE LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) <= 1 -- 숫자가 0개 또는 1개인
)
SELECT
ROUND(COUNT(*) / (SELECT COUNT(*) FROM `proven-serenity-439401-f2.modulabs_project.data`) * 100, 2) AS
FROM
`proven-serenity-439401-f2.modulabs_project.data`
WHERE
StockCode IN (SELECT StockCode FROM StockCodesWithFewNumbers);
```

쿼리 결과



• 제품과 관련되지 않은 거래 기록을 제거하기

```
DELETE FROM proven-serenity-439401-f2.modulabs_project.data_re
WHERE StockCode IN (
    SELECT DISTINCT StockCode
    FROM (
        SELECT StockCode
    FROM proven-serenity-439401-f2.modulabs_project.data_re
```

```
WHERE StockCode IN ('POST', 'BANK CHARGES', 'PADS', 'DOT', 'CARRIAGE', 'D', 'M', 'CRUK') -- 제품과 관
) AS non_product_codes
);
```



처음부터 데이터셋 새로 올려서 다시 수행했는데로 아래와 같이 나옴



Description 살펴보기

• 고유한 Description 별 출현 빈도를 계산하고 상위 30개를 출력하기

```
SELECT Description, COUNT(*) AS description_cnt
FROM proven-serenity-439401-f2.modulabs_project.data_re
GROUP BY Description
ORDER BY description_cnt DESC
LIMIT 30;
```

| 쿼리 | 결과 | | | | | |
|------|------------------|-------------|-------------|---------|-----|--------|
| 작업 정 | 보 결과 | 차트 | JSON | 실행 세부 | ·정보 | 실행 그래프 |
| 행 // | Description ▼ | // | description | n_cnt ▼ | | |
| 1 | WHITE HANGING | HEART T-LIG | | 2058 | | |
| 2 | REGENCY CAKES | TAND 3 TIER | | 1894 | | |
| 3 | JUMBO BAG RED | RETROSPOT | | 1659 | | |
| 4 | PARTY BUNTING | | | 1409 | | |
| 5 | ASSORTED COLO | UR BIRD ORN | | 1405 | | |
| 6 | LUNCH BAG RED | RETROSPOT | | 1345 | | |
| 7 | SET OF 3 CAKE TI | NS PANTRY | | 1224 | | |
| 8 | LUNCH BAG BLA | CK SKULL. | | 1099 | | |
| 더보기 | | | | | | |

• 서비스 관련 정보를 포함하는 행들을 제거하기

```
DELETE
FROM proven-serenity-439401-f2.modulabs_project.data_re
WHERE
Description LIKE '%Next Day Carriage%' OR
Description LIKE '%High Resolution Image%';
```

- 캡처 못해서 , 다시 실행 시 삭제 행이 '0'으로 나옴
- 대소문자를 혼합하고 있는 데이터를 대문자로 표준화 하기

```
CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.data_re AS
SELECT

* EXCEPT (Description),
UPPER(Description) AS Description -- Description을 대문자로 변환
FROM proven-serenity-439401-f2.modulabs_project.data_re;
```

```
198 -- 대소문자를 혼합하고 있는 데이터를 대문자로 표준화
199 CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.data_re AS
200 SELECT
201 * EXCEPT (Description),
102 UPPER(Description) AS Description -- Description을 대문자로 변환
FROM proven-serenity-439401-f2.modulabs_project.data_re;

라 결과

작업 정보 결과 실행 세부정보 실행 그래프

① 이 문으로 이름이 data_re인 테이블이 교체되었습니다.
```

UnitPrice 살펴보기

• UnitPrice 의 최솟값, 최댓값, 평균을 구하기

```
SELECT

MIN(UnitPrice) AS min_price,

MAX(UnitPrice) AS max_price,

AVG(UnitPrice) AS avg_price

FROM proven-serenity-439401-f2.modulabs_project.data_re;
```

쿼리 결과



• 단가가 0원인 거래의 개수, 구매 수량(Quantity)의 최솟값, 최댓값, 평균 구하기

```
SELECT

COUNT(*) AS cnt_quantity, -- 단가가 0원인 거래의 개수
MIN(Quantity) AS min_quantity, -- 구매 수량의 최솟값
MAX(Quantity) AS max_quantity, -- 구매 수량의 최댓값
AVG(Quantity) AS avg_quantity -- 구매 수량의 평균값
FROM proven-serenity-439401-f2.modulabs_project.data_re
WHERE UnitPrice = 0; -- 단가가 0원인 거래 필터링
```

| 쿼리 | 결과 | | | | | |
|------|-------|------------|--------|-----------|----------------|-------------------|
| 작업 정 | 경보 | 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
| 행 | cnt_c | quantity 🔻 | min_qı | uantity ▼ | max_quantity ▼ | avg_quantity ▼ |
| 1 | | 33 | | 1 | 12540 | 420.5151515151515 |

• UnitPrice = 0 를 제거하고 일관된 데이터셋을 유지하기

CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.data_re AS
SELECT *
FROM proven-serenity-439401-f2.modulabs_project.data_re
WHERE UnitPrice != 0;

쿼리 결과





11-7. RFM 스코어

Recency

쿼리 결과

• InvoiceDate 컬럼을 연월일 자료형으로 변경하기

```
SELECT
DATE(InvoiceDate) AS InvoiceDay, -- InvoiceDate에서 연월일만 추출
*
FROM proven-serenity-439401-f2.modulabs_project.data_re;
```

▲ 결과 저장 ▼

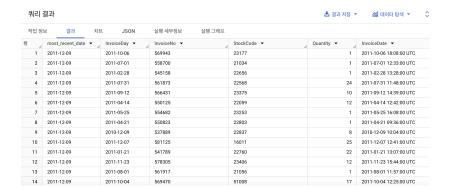
₩ 데이터 탐색 ▼

| 14-1 | 2-1 | | | | □ 일파 시장 · | 10대 대이다 함색 * |
|------|--------------|-------------|----------------|------------|-------------------------|--------------|
| 작업 정 | j보 결과 | 차트 JSON | 실행 세부정보 실행 그래프 | | | |
| 행 | InvoiceDay ▼ | InvoiceNo ▼ | StockCode ▼ | Quantity • | InvoiceDate ▼ | UnitPrice ▼ |
| 1 | 2011-11-03 | 574301 | 85049E | 12 | 2011-11-03 16:15:00 UTC | 1.25 |
| 2 | 2011-11-03 | 574301 | 85049A | 12 | 2011-11-03 16:15:00 UTC | 1.25 |
| 3 | 2011-11-03 | 574301 | 22910 | 6 | 2011-11-03 16:15:00 UTC | 2.95 |
| 4 | 2011-11-03 | 574301 | 23514 | 6 | 2011-11-03 16:15:00 UTC | 2.08 |
| 5 | 2011-11-03 | 574301 | 22086 | 6 | 2011-11-03 16:15:00 UTC | 2.95 |
| 6 | 2011-11-03 | 574301 | 20749 | 4 | 2011-11-03 16:15:00 UTC | 7.95 |
| 7 | 2011-11-03 | 574301 | 23240 | 6 | 2011-11-03 16:15:00 UTC | 4.15 |
| 8 | 2011-11-03 | 574301 | 20971 | 12 | 2011-11-03 16:15:00 UTC | 1.25 |
| 9 | 2011-11-03 | 574301 | 22751 | 4 | 2011-11-03 16:15:00 UTC | 3.75 |
| 10 | 2011-11-03 | 574301 | 22144 | 6 | 2011-11-03 16:15:00 UTC | 2.1 |
| 11 | 2011-11-03 | 574301 | 23511 | 6 | 2011-11-03 16:15:00 UTC | 2.08 |
| 12 | 2011-11-03 | 574301 | 22077 | 12 | 2011-11-03 16:15:00 UTC | 1.95 |

• 가장 최근 구매 일자를 MAX() 함수로 찾아보기

SELECT

MAX(DATE(InvoiceDate)) OVER() AS most_recent_date, -- 가장 최근 구매 일자를 윈도우 함수로 계산
DATE(InvoiceDate) AS InvoiceDay, -- 연월일로 변환한 구매 일자
*
FROM proven-serenity-439401-f2.modulabs_project.data_re;



• 유저 별로 가장 큰 InvoiceDay를 찾아서 가장 최근 구매일로 저장하기

```
SELECT
CustomerID,
MAX(DATE(InvoiceDate)) AS most_recent_date -- 유저별 가장 최근 구매일을 most_recent_date로 저장
```

```
FROM proven-serenity-439401-f2.modulabs_project.data_re
GROUP BY CustomerID; -- 유저별로 그룹화
```

| 작업 정보 | 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
|-------|--------------|------------|---------|---------|--------|
| 행 // | CustomerID ▼ | most_recen | nt_date | | |
| 1 | 12544 | 2011-11-10 | | | |
| 2 | 13568 | 2011-06-19 | | | |
| 3 | 13824 | 2011-11-07 | | | |
| 4 | 14080 | 2011-11-07 | | | |
| 5 | 14336 | 2011-11-23 | | | |
| 6 | 14592 | 2011-11-04 | | | |
| 7 | 15104 | 2011-06-26 | | | |
| 8 | 15360 | 2011-10-31 | | | |
| 9 | 15872 | 2011-11-25 | | | |
| 10 | 16128 | 2011-11-22 | | | |
| 11 | 16384 | 2011-09-11 | | | |
| 12 | 17152 | 2011-05-29 | | | |

• 가장 최근 일자(most_recent_date)와 유저별 마지막 구매일(InvoiceDay)간의 차이를 계산하기

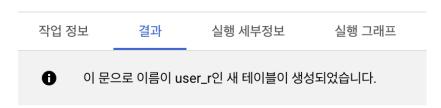
```
SELECT
CustomerID,
EXTRACT(DAY FROM MAX(InvoiceDay) OVER () - InvoiceDay) AS recency -- 전체에서 가장 최근 일자와 유저별 구매일
FROM (
SELECT
CustomerID,
MAX(DATE(InvoiceDate)) AS InvoiceDay -- 유저별 마지막 구매일
FROM proven-serenity-439401-f2.modulabs_project.data_re
GROUP BY CustomerID
);
```

| 작업 정보 | 럳 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
|-------|--------------|---------|----------|---------|--------|
| 행 // | CustomerID ▼ | recency | ▼ | | |
| 1 | 13824 | | 32 | | |
| 2 | 12551 | | 357 | | |
| 3 | 14345 | | 38 | | |
| 4 | 17418 | | 42 | | |
| 5 | 13579 | | 14 | | |
| 6 | 15380 | | 8 | | |
| 7 | 16150 | | 38 | | |
| 8 | 16919 | | 156 | | |
| 9 | 14364 | | 108 | | |
| 10 | 15667 | | 39 | | |
| 11 | 17736 | | 9 | | |

• 최종 데이터 셋에 필요한 데이터들을 각각 정제해서 이어붙이고 지금까지의 결과를 user_r 이라는 이름의 테이블로 저장하기

```
CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.user_r AS
SELECT
CustomerID,
EXTRACT(DAY FROM MAX(InvoiceDay) OVER () - InvoiceDay) AS recency -- 가장 최근 일자와 유저별 구매일 차이 계상
FROM (
SELECT
CustomerID,
MAX(DATE(InvoiceDate)) AS InvoiceDay -- 유저별 마지막 구매일
FROM proven-serenity-439401-f2.modulabs_project.data_re
GROUP BY CustomerID
);
```

쿼리 결과





| 스키마 | 세부정보 | 미리보기 |
|------|------------|---------|
| 행 // | CustomerID | recency |
| 1 | 12985 | 0 |
| 2 | 18102 | 0 |
| 3 | 14051 | 0 |
| 4 | 17754 | 0 |
| 5 | 14446 | 0 |
| 6 | 13426 | 0 |
| 7 | 16446 | 0 |
| 8 | 17001 | 0 |
| 9 | 15804 | 0 |
| 10 | 14422 | 0 |
| 11 | 14441 | 0 |
| 12 | 16954 | 0 |
| 13 | 13777 | 0 |
| 14 | 12518 | 0 |
| 15 | 12433 | 0 |

Frequency

• 고객마다 고유한 InvoiceNo의 수를 세어보기

```
SELECT
CustomerID,
COUNT(DISTINCT InvoiceNo) AS purchase_cnt -- 고객별 고유한 InvoiceNo 수(거래 건수)
FROM proven-serenity-439401-f2.modulabs_project.data_re
GROUP BY CustomerID;
```

| 작업 정! | 보 결과 | 차트 JSC | ON |
|-------|--------------|--------------|----------|
| 행 | CustomerID ▼ | purchase_cnt | - |
| 1 | 12544 | | 2 |
| 2 | 13568 | | 1 |
| 3 | 13824 | | 5 |
| 4 | 14080 | | 1 |
| 5 | 14336 | | 4 |
| 6 | 14592 | | 3 |
| 7 | 15104 | | 3 |
| 8 | 15360 | | 1 |
| 9 | 15872 | | 2 |
| 10 | 16128 | | 5 |
| 11 | 16384 | | 2 |
| 12 | 17152 | | 4 |
| 13 | 17408 | | 1 |

• 각 고객 별로 구매한 아이템의 총 수량 더하기

```
SELECT
CustomerID,
SUM(Quantity) AS item_cnt -- 고객별로 구매한 아이템의 총 수량 합계
FROM proven-serenity-439401-f2.modulabs_project.data_re
GROUP BY CustomerID;
```

쿼리 결과

| 작업 정보 | 결과 | 차트 | JSON | 실행 세부정보 | 실행 그래프 |
|-------|--------------|----------|----------|---------|--------|
| 행 | CustomerID ▼ | item_cnt | ~ | | |
| 1″ | 12544 ~ | | 130 ″ | | |
| 2 | 13568 | | 66 | | |
| 3 | 13824 | | 768 | | |
| 4 | 14080 | | 48 | | |
| 5 | 14336 | | 1759 | | |
| 6 | 14592 | | 407 | | |
| 7 | 15104 | | 633 | | |
| 8 | 15360 | | 223 | | |
| 9 | 15872 | | 187 | | |
| 10 | 16128 | | 988 | | |
| 11 | 16384 | | 260 | | |
| 12 | 17152 | | 477 | | |
| 13 | 17408 | | 3 | | |

• 전체 거래 건수 계산와 구매한 아이템의 총 수량 계산의 결과를 합쳐서 user_rf 라는 이름의 테이블에 저장하기

```
CREATE OR REPLACE TABLE `proven-serenity-439401-f2.modulabs_project.user_rf` AS
-- (1) 전체 거래 건수 계산
WITH purchase_cnt AS (
 SELECT
   CustomerID,
   COUNT(DISTINCT InvoiceNo) AS purchase_cnt -- 고객별 고유한 InvoiceNo 수(거래 건수)
  FROM `proven-serenity-439401-f2.modulabs_project.data_re`
 GROUP BY CustomerID
),
-- (2) 구매한 아이템 총 수량 계산
item_cnt AS (
 SELECT
   CustomerID,
   SUM(Quantity) AS item_cnt -- 고객별로 구매한 아이템의 총 수량 합계
 FROM `proven-serenity-439401-f2.modulabs_project.data_re`
 GROUP BY CustomerID
-- 기존의 user_r에 (1)과 (2)를 통합
SELECT
 pc.CustomerID,
 pc.purchase_cnt,
 ic.item_cnt,
 ur.recency
FROM purchase_cnt AS pc
JOIN item_cnt AS ic
 ON pc.CustomerID = ic.CustomerID
JOIN `proven-serenity-439401-f2.modulabs_project.user_r` AS ur
 ON pc.CustomerID = ur.CustomerID;
```

작업 정보 결과 실행 세부정보 실행 그래프 이 문으로 이름이 user_rf인 새 테이블이 생성되었습니다.

| ⊞ us | er_rf (| Q.쿼리 ▼ +2 | 공유 □복사 | ▋스냅샷 |
|-------------|------------|--------------|--------------|---------|
| 스키마 | 세부정보 | 미리보기 | 테이블 탐색기 미리보기 | |
| 행 // | CustomerID | purchase_cnt | item_cnt | recency |
| 1 | 12713 | 1 | 505 | 0 |
| 2 | 12792 | 1 | 215 | 256 |
| 3 | 18010 | 1 | 60 | 256 |
| 4 | 15083 | 1 | 38 | 256 |
| 5 | 15520 | 1 | 314 | 1 |
| 6 | 14569 | 1 | 79 | 1 |
| 7 | 13298 | 1 | 96 | 1 |
| 8 | 13436 | 1 | 76 | 1 |
| 9 | 14476 | 1 | 110 | 257 |
| 10 | 13357 | 1 | 321 | 257 |
| 11 | 14204 | 1 | 72 | 2 |
| 12 | 15195 | 1 | 1404 | 2 |
| 13 | 15471 | 1 | 256 | 2 |
| 14 | 12442 | 1 | 181 | 3 |
| 15 | 12650 | 1 | 250 | 3 |

Monetary

• 고객별 총 지출액 계산 (소수점 첫째 자리에서 반올림)

```
SELECT
CustomerID,
ROUND(SUM(UnitPrice * Quantity), 1) AS user_total -- 고객별 총 지출액, 소수점 첫째 자리에서 반올림
FROM proven-serenity-439401-f2.modulabs_project.data_re
GROUP BY CustomerID;
```

| 작업 정보 | 년 결 <u>과</u> | 차트 JSON | 실행 세부정보 | 실행 그래프 |
|-------|--------------|--------------|---------|--------|
| H // | CustomerID ▼ | user_total ▼ | | |
| 1 | 12544 | 299.7 | | |
| 2 | 13568 | 187.1 | | |
| 3 | 13824 | 1698.9 | | |
| 4 | 14080 | 45.6 | | |
| 5 | 14336 | 1614.9 | | |
| 6 | 14592 | 557.9 | | |
| 7 | 15104 | 968.6 | | |
| 8 | 15360 | 427.9 | | |
| 9 | 15872 | 316.2 | | |
| 10 | 16128 | 1880.2 | | |

• 고객별 평균 거래 금액 계산

○ 고객별 평균 거래 금액을 구하기 위해 1) data 테이블을 user_rf 테이블과 조인(LEFT JOIN) 한 후, 2) purchase_cnt 로 나누어서 3) user_rfm 테이블로 저장하기

```
CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.user_rfm AS
 rf.CustomerID AS CustomerID,
 rf.purchase_cnt,
 rf.item_cnt,
 rf.recency,
 ut.user_total,
 ROUND(ut.user_total / rf.purchase_cnt, 1) AS user_average -- 평균 거래 금액 계산
FROM proven-serenity-439401-f2.modulabs_project.user_rf rf
LEFT JOIN (
 -- 고객 별 총 지출액
 SELECT
   CustomerID,
   ROUND(SUM(UnitPrice * Quantity), 1) AS user_total -- 고객별 총 지출액 계산, 소수점 첫째 자리 반올림
 FROM proven-serenity-439401-f2.modulabs_project.data_re
 GROUP BY CustomerID
) ut
ON rf.CustomerID = ut.CustomerID;
```

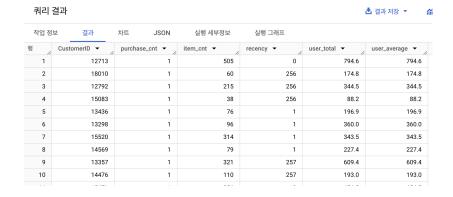




RFM 통합 테이블 출력하기

• 최종 user_rfm 테이블을 출력하기

```
SELECT *
FROM proven-serenity-439401-f2.modulabs_project.user_rfm
LIMIT 100;
```



11-8. 추가 Feature 추출

1. 구매하는 제품의 다양성

1) 고객 별로 구매한 상품들의 고유한 수를 계산하기 2)
 user_rfm 테이블과 결과를 합치기 3)
 user_data 라는 이름의 테이블에 저장하기

```
CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.user_data AS
WITH unique_products AS (
SELECT
CustomerID,
COUNT(DISTINCT StockCode) AS unique_products -- 고유한 제품 수 계산
```

```
FROM proven-serenity-439401-f2.modulabs_project.data_re
    GROUP BY CustomerID
)
SELECT
    ur.*,
    up.unique_products -- unique_products 값을 추가
FROM proven-serenity-439401-f2.modulabs_project.user_rfm AS ur
JOIN unique_products AS up
ON ur.CustomerID = up.CustomerID;
```

[결과 이미지를 넣어주세요]



2. 평균 구매 주기

- 고객들의 쇼핑 패턴을 이해하는 것을 목표 (고객 별 재방문 주기 살펴보기)
 - 평균 구매 소요 일수를 계산하고, 그 결과를 user_data 에 통합

```
CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.user_data AS
WITH purchase_intervals AS (
  -- (2) 고객 별 구매와 구매 사이의 평균 소요 일수
 SELECT
   CustomerID,
   CASE WHEN ROUND(AVG(interval_), 2) IS NULL THEN 0 ELSE ROUND(AVG(interval_), 2) END AS average_interval
   -- (1) 구매와 구매 사이에 소요된 일수
   SELECT
     CustomerID,
     DATE_DIFF(InvoiceDate, LAG(InvoiceDate) OVER (PARTITION BY CustomerID ORDER BY InvoiceDate), DAY) AS
     proven-serenity-439401-f2.modulabs_project.data_re
   WHERE CustomerID IS NOT NULL
 GROUP BY CustomerID
SELECT u.*, pi.* EXCEPT (CustomerID)
FROM proven-serenity-439401-f2.modulabs_project.user_data AS u
LEFT JOIN purchase_intervals AS pi
ON u.CustomerID = pi.CustomerID;
```

쿼리 결과 작업 정보 결과 실행 세부정보 실행 그래프 ① 문으로 이름이 user_data인 테이블이 교체되었습니다.



3. 구매 취소 경향성

- 고객의 취소 패턴 파악하기
 - 1) 취소 빈도(cancel_frequency) : 고객 별로 취소한 거래의 총 횟수
 - 2) 취소 비율(cancel_rate): 각 고객이 한 모든 거래 중에서 취소를 한 거래의 비율
 - 취소 빈도와 취소 비율을 계산하고 그 결과를 user_data 에 통합하기 (취소 비율은 소수점 두번째 자리)

```
CREATE OR REPLACE TABLE proven-serenity-439401-f2.modulabs_project.user_data AS
WITH TransactionInfo AS (
 SELECT
   CustomerID,
   COUNT(*) AS total_transactions, -- 총 거래 횟수
   SUM(CASE WHEN InvoiceNo LIKE 'C%' THEN 1 ELSE 0 END) AS cancel_frequency -- 취소된 거래 횟수
 FROM proven-serenity-439401-f2.modulabs_project.data_re
 WHERE CustomerID IS NOT NULL
 GROUP BY CustomerID
)
SELECT
 u.*,
 t.* EXCEPT(CustomerID),
 ROUND(t.cancel_frequency / t.total_transactions, 2) AS cancel_rate -- 취소 비율 계산 (소수점 두번째 자리)
FROM proven-serenity-439401-f2.modulabs_project.user_data AS u
LEFT JOIN TransactionInfo AS t
ON u.CustomerID = t.CustomerID;
```



• 다양한 컬럼들을 활용하여 고객의 구매 패턴과 선호도를 보다 심층적으로 이해할 수 있도록 최종적으로 user_data 를 출력하기

SELECT * FROM proven-serenity-439401-f2.modulabs_project.user_data LIMIT 100;



회고

이번 프로젝트를 통해 고객 세그먼테이션을 위한 **RFM 분석**과 추가적인 **고객 구매 패턴** 분석을 진행했습니다. RFM 분석 외에도 **구매 제품의 다양성, 평균 구매 주기, 취소 패턴** 등의 추가 특성을 추출함으로써 고객을 더 깊이 이해할 수 있었습니다. 각 특성을 바탕으로 **클러스터링 알고리즘**을 적용해 고객을 그룹화할 수 있는 기반을 마련했고, SQL을 활용한 데이터 처리와 분석이 실제로 어떻게 마케팅 전략에 적용될 수 있는지 실습할 수 있었습니다.