

SQL_11. 고객을 세그먼테이션하자![프로젝트]

-- 데이터 살펴보기

```
SELECT *  
FROM `modulabs_project.data`  
LIMIT 10
```

작업 정보	결과	차트	JSON	실행 세부정보	실행 그래프
행	InvoiceNo	StockCode	Description	Quantity	
1	536365	85123A	WHITE HANGING HEART T-LIG...	6	
2	536365	71053	WHITE METAL LANTERN	6	
3	536365	84406B	CREAM CUPID HEARTS COAT H...	8	
4	536365	84029G	KNITTED UNION FLAG HOT WA...	6	
5	536365	84029E	RED WOOLLY HOTTIE WHITE H...	6	
6	536365	22752	SET 7 BABUSHKA NESTING BO...	2	
7	536365	21730	GLASS STAR FROSTED T-LIGHT...	6	
8	536366	22633	HAND WARMER UNION JACK	6	

-- 데이터 행 구성 보기

```
SELECT COUNT(*) AS row_count  
FROM `modulabs-project-465302.modulabs_project.data`
```

행	row_count
1	541909

-- 데이터 수 세기

```
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 `modulabs-project-465302.modulabs_project.data`

```

행	count_InvoiceNo	count_StockCode	count_Description	count_Quantity	count_InvoiceDate	count_UnitPrice
1	541909	541909	540455	541909	541909	541

```

-- 데이터 전처리 _ 결측치 제거
SELECT
  'InvoiceNo' AS column_name,
  ROUND(SUM(CASE WHEN InvoiceNo IS NULL THEN 1 ELSE 0 END) / COUNT(*)) AS missing_percentage
FROM `modulabs-project-465302.modulabs_project.data`

```

행	column_name	missing_percentage
1	InvoiceNo	0.0

```

-- 결측치 알아보기
-- 다른 컬럼에도 동일하게 반영한 후, UNION ALL로 연결
SELECT
  'InvoiceNo' AS column_name,
  ROUND(SUM(CASE WHEN InvoiceNo IS NULL THEN 1 ELSE 0 END) / COUNT(*)) AS missing_percentage
FROM `modulabs-project-465302.modulabs_project.data`

```

UNION ALL

```

SELECT
  'StockCode',

```

```
ROUND(SUM(CASE WHEN StockCode IS NULL THEN 1 ELSE 0 END) / COUNT(*)  
FROM `modulabs-project-465302.modulabs_project.data`
```

UNION ALL

```
SELECT  
  'Description',  
  ROUND(SUM(CASE WHEN Description IS NULL THEN 1 ELSE 0 END) / COUNT(*)  
FROM `modulabs-project-465302.modulabs_project.data`
```

UNION ALL

```
SELECT  
  'Quantity',  
  ROUND(SUM(CASE WHEN Quantity IS NULL THEN 1 ELSE 0 END) / COUNT(*) *  
FROM `modulabs-project-465302.modulabs_project.data`
```

UNION ALL

```
SELECT  
  'InvoiceDate',  
  ROUND(SUM(CASE WHEN InvoiceDate IS NULL THEN 1 ELSE 0 END) / COUNT(*)  
FROM `modulabs-project-465302.modulabs_project.data`
```

UNION ALL

```
SELECT  
  'UnitPrice',  
  ROUND(SUM(CASE WHEN UnitPrice IS NULL THEN 1 ELSE 0 END) / COUNT(*)  
FROM `modulabs-project-465302.modulabs_project.data`
```

UNION ALL

```
SELECT  
  'CustomerID',  
  ROUND(SUM(CASE WHEN CustomerID IS NULL THEN 1 ELSE 0 END) / COUNT(*)
```

```
FROM `modulabs-project-465302.modulabs_project.data`
```

```
UNION ALL
```

```
SELECT
```

```
  'Country',
```

```
  ROUND(SUM(CASE WHEN Country IS NULL THEN 1 ELSE 0 END) / COUNT(*) *
```

```
FROM `modulabs-project-465302.modulabs_project.data`
```

행	column_name ▼	missing_percenta...
1	InvoiceNo	0.0
2	InvoiceDate	0.0
3	Country	0.0
4	StockCode	0.0
5	Quantity	0.0
6	Description	0.27
7	UnitPrice	0.0
8	CustomerID	24.93

```
-- 같은 제품(StockCode)이 항상 같은 상세 설명(Description)을 가지고 있지 않다는 데이
```

```
-- StockCode = '85123A'의 Description을 추출하는 쿼리문을 작성
```

```
SELECT Description
```

```
FROM `modulabs-project-465302.modulabs_project.data`
```

```
WHERE StockCode = '85123A'
```

행	Description ▼	
1	WHITE HANGING HEART T-LIGHT HOLDER	
2	WHITE HANGING HEART T-LIGHT HOLDER	
3	WHITE HANGING HEART T-LIGHT HOLDER	
4	WHITE HANGING HEART T-LIGHT HOLDER	
5	WHITE HANGING HEART T-LIGHT HOLDER	
6	WHITE HANGING HEART T-LIGHT HOLDER	
7	WHITE HANGING HEART T-LIGHT HOLDER	
8	WHITE HANGING HEART T-LIGHT HOLDER	

-- 결측치 처리

```
DELETE FROM `modulabs-project-465302.modulabs_project.data`
WHERE InvoiceNo IS NULL
  OR StockCode IS NULL
  OR Description IS NULL
  OR Quantity IS NULL
  OR InvoiceDate IS NULL
  OR UnitPrice IS NULL
  OR CustomerID IS NULL
  OR Country IS NULL
```

i 이 문으로 data의 행 135,080개가 삭제되었습니다.

-- 데이터 전처리_ 중복값 확인

-- 중복된 행의 수를 세기

```
SELECT
  COUNT(*) AS duplicate_count
FROM (
  SELECT
    InvoiceNo,
```

```

StockCode,
Description,
Quantity,
InvoiceDate,
UnitPrice,
CustomerID,
Country,
COUNT(*) AS cnt
FROM `modulabs-project-465302.modulabs_project.data`
GROUP BY
    InvoiceNo, StockCode, Description, Quantity, InvoiceDate, UnitPrice, CustomerID, Country
HAVING COUNT(*) > 1
)

```

행	duplicate_count
1	4837

```

-- 중복값 처리
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.data`
SELECT DISTINCT *
FROM `modulabs-project-465302.modulabs_project.data`

```

i 이 문으로 이름이 data인 테이블이 교체되었습니다.

```

-- 중복값 처리 이후 남은 행의 개수
SELECT COUNT(*) AS row_count
FROM `modulabs-project-465302.modulabs_project.data`

```

행	row_count
1	401604

```
-- 데이터 전처리_오류값 처리
-- InvoiceNo 살펴보기

-- 고유(unique)한 InvoiceNo의 개수를 출력
SELECT
  COUNT(DISTINCT InvoiceNo) AS unique_invoice_count
FROM `modulabs-project-465302.modulabs_project.data`
```

행	unique_invoice_c...
1	22190

```
-- 고유한 InvoiceNo를 100개를 출력
SELECT DISTINCT InvoiceNo
FROM `modulabs-project-465302.modulabs_project.data`
LIMIT 100
```

행	InvoiceNo ▼
1	541431
2	C541433
3	537626
4	542237
5	549222
6	556201
7	562032
8	573511

```
-- InvoiceNo가 'C'로 시작하는 행을 필터링
SELECT *
FROM `modulabs-project-465302.modulabs_project.data`
```

```
WHERE InvoiceNo LIKE 'C%'
LIMIT 100;
```

행	InvoiceNo	StockCode	Description	Quantity
1	C541433	23166	MEDIUM CERAMIC TOP STORA...	-74215
2	C545329	M	Manual	-1
3	C545329	M	Manual	-1
4	C545330	M	Manual	-1
5	C547388	84050	PINK HEART SHAPE EGG FRYIN...	-12
6	C547388	22645	CERAMIC HEART FAIRY CAKE ...	-12
7	C547388	22784	LANTERN CREAM GAZEBO	-3
8	C547388	21914	BLUE HARMONICA IN BOX	-12

```
-- 구매 건 상태가 Canceled 인 데이터의 비율(%)
SELECT ROUND(
  SUM(CASE WHEN InvoiceNo LIKE 'C%' THEN 1 ELSE 0 END)
  / COUNT(*) * 100,
  1
) AS canceled_rate_percentage
FROM `modulabs-project-465302.modulabs_project.data`
```

행	canceled_rate_pe...
1	2.2

```
-- StockCode 살펴보기
```

```
-- 고유한 StockCode의 개수를 출력
SELECT COUNT(DISTINCT StockCode) AS unique_stockcode_count
FROM `modulabs-project-465302.modulabs_project.data`
```

행	unique_stockcod...
1	3684


```
-- 어떤 제품이 가장 많이 판매되었는지 보기 위하여 StockCode 별 등장 빈도를 출력(상위 10개)
SELECT
  StockCode,
  COUNT(*) AS frequency
FROM `modulabs-project-465302.modulabs_project.data`
GROUP BY StockCode
ORDER BY frequency DESC
LIMIT 10
```

작업 정보	결과	차트	JSON	실행 세부정보
행	StockCode ▼	frequency ▼		
1	85123A	2065		
2	22423	1894		
3	85099B	1659		
4	47566	1409		
5	84879	1405		
6	20725	1346		
7	22720	1224		
8	POST	1196		

```
--이상치들이 몇 개나 있는지 확인하기 위하여 StockCode의 문자열 내 숫자의 길이 출력
WITH UniqueStockCodes AS (
  SELECT DISTINCT StockCode
  FROM `modulabs-project-465302.modulabs_project.data`
)
SELECT
  LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS
  COUNT(*) AS stock_cnt
FROM UniqueStockCodes
GROUP BY number_count
ORDER BY stock_cnt DESC
```

행	number_count	stock_cnt
1	5	3676
2	0	7
3	1	1

-- 출력 결과를 보면, 8개를 제외하곤 StockCode에 5개의 숫자들이 포함되어 있는 것을 알
 -- 숫자가 0개인 코드는 7개, 숫자가 1개인 코드는 1개

-- 숫자가 0~1개인 값들에는 어떤 코드들이 들어가 있는지를 확인

```
SELECT DISTINCT StockCode, number_count
FROM (
  SELECT
    StockCode,
    LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS
  FROM `modulabs-project-465302.modulabs_project.data`
)
WHERE number_count <= 1
```

행	StockCode	number_count
1	POST	0
2	M	0
3	C2	1
4	D	0
5	BANK CHARGES	0
6	PADS	0
7	DOT	0
8	CRUK	0

-- 데이터 수는 전체 데이터 수 대비 몇 퍼센트?

```
SELECT
  ROUND(
    COUNTIF(
```

```

    LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) <
    ) / COUNT(*) * 100,
    2
    ) AS special_stockcode_percentage
FROM `modulabs-project-465302.modulabs_project.data`

```

행	special_stockcod...
1	0.48

```

-- 제품과 관련되지 않은 거래 기록을 제거
DELETE FROM `modulabs-project-465302.modulabs_project.data`
WHERE StockCode IN (
    SELECT DISTINCT StockCode
    FROM (
        SELECT
            StockCode,
            LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS
        FROM `modulabs-project-465302.modulabs_project.data`
    )
    WHERE number_count <= 1
)

```

 이 문으로 data의 행 1,915개가 삭제되었습니다.

```

-- Description 살펴보기
-- 고유한 Description 별 출현 빈도를 계산하고 상위 30개를 출력
SELECT
    Description,
    COUNT(*) AS frequency
FROM `modulabs-project-465302.modulabs_project.data`
GROUP BY Description

```

ORDER BY frequency DESC
LIMIT 30

행	Description ▼	frequency ▼
1	WHITE HANGING HEART T-LIG...	2058
2	REGENCY CAKESTAND 3 TIER	1894
3	JUMBO BAG RED RETROSPOT	1659
4	PARTY BUNTING	1409
5	ASSORTED COLOUR BIRD ORN...	1405
6	LUNCH BAG RED RETROSPOT	1345
7	SET OF 3 CAKE TINS PANTRY D...	1224
8	LUNCH BAG BLACK SKULL.	1099

-- 대소문자가 혼합된 Description이 있는지 확인
SELECT DISTINCT Description
FROM `modulabs-project-465302.modulabs_project.data`
WHERE REGEXP_CONTAINS(Description, r'[a-z]');

행	Description ▼
1	BAG 250g SWIRLY MARBLES
2	3 TRADITIONAL BISCUIT CUTTE...
3	BAG 125g SWIRLY MARBLES
4	POLYESTER FILLER PAD 30CMx...
5	BAG 500g SWIRLY MARBLES
6	POLYESTER FILLER PAD 45x45...
7	POLYESTER FILLER PAD 40x40...
8	ESSENTIAL BALM 3.5g TIN IN E...

-- 서비스 관련 정보를 포함하는 행들을 제거
DELETE

```

FROM `modulabs-project-465302.modulabs_project.data`
WHERE
  UPPER(Description) LIKE '%POSTAGE%' OR
  UPPER(Description) LIKE '%CARRIAGE%' OR
  UPPER(Description) LIKE '%BANK CHARGES%' OR
  UPPER(Description) LIKE '%ADJUST%' OR
  UPPER(Description) LIKE '%MANUAL%' OR
  UPPER(Description) LIKE '%CHECK%' OR
  UPPER(Description) LIKE '%SAMPLES%' OR
  UPPER(Description) LIKE '%GIFT%'


```

 이 문으로 data의 행 4,202개가 삭제되었습니다.

```

-- 대소문자를 혼합하고 있는 데이터를 대문자로 표준화
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.data`
SELECT
  * EXCEPT (Description),
  UPPER(Description) AS Description
FROM `modulabs-project-465302.modulabs_project.data`

```

 이 문으로 이름이 data인 테이블이 교체되었습니다.

```

-- UnitPrice 살펴보기
-- UnitPrice에서 이상치 찾기

-- UnitPrice의 최솟값, 최댓값, 평균
SELECT
  MIN(UnitPrice) AS min_price,
  MAX(UnitPrice) AS max_price,
  ROUND(AVG(UnitPrice), 2) AS avg_price
FROM `modulabs-project-465302.modulabs_project.data`

```

행	min_price	max_price	avg_price
1	0.0	649.5	2.9

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

```
SELECT
  COUNT(*) AS cnt_quantity,
  MIN(Quantity) AS min_quantity,
  MAX(Quantity) AS max_quantity,
  ROUND(AVG(Quantity), 2) AS avg_quantity
FROM `modulabs-project-465302.modulabs_project.data`
WHERE UnitPrice = 0
```

행	cnt_quantity	min_quantity	max_quantity	avg_quantity
1	33	1	12540	420.52

-- 이 데이터(UnitPrice = 0)를 제거하고 일관된 데이터셋을 유지

```
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.data`
SELECT *
FROM `modulabs-project-465302.modulabs_project.data`
WHERE UnitPrice > 0
```

i 이 문으로 이름이 data인 테이블이 교체되었습니다.

-- RFM 스코어

-- Recency

-- InvoiceDate 컬럼을 연월일 자료형으로 변경

```
SELECT
  DATE(InvoiceDate) AS InvoiceDay,
```

*

```
FROM `modulabs-project-465302.modulabs_project.data`
```

행	InvoiceDay	InvoiceNo	StockCode	Quantity	InvoiceDate
1	2011-01-18	541431	23166	74215	2011-01-18
2	2011-01-18	C541433	23166	-74215	2011-01-18
3	2010-12-07	537626	22727	4	2010-12-07
4	2010-12-07	537626	22212	6	2010-12-07
5	2010-12-07	537626	22195	12	2010-12-07
6	2010-12-07	537626	21171	12	2010-12-07
7	2010-12-07	537626	22726	4	2010-12-07
8	2010-12-07	537626	22774	12	2010-12-07

```
-- 가장 최근 구매 일자를 MAX() 함수로 찾기
```

```
SELECT
```

```
MAX(DATE(InvoiceDate)) OVER () AS most_recent_date,
```

```
DATE(InvoiceDate) AS InvoiceDay,
```

*

```
FROM `modulabs-project-465302.modulabs_project.data`
```

행	most_recent_date	InvoiceDay	InvoiceNo	StockCode	Quantity
1	2011-12-09	2011-08-02	562046	84792	
2	2011-12-09	2011-11-24	578459	22338	
3	2011-12-09	2011-07-22	560991	16219	
4	2011-12-09	2011-09-15	566773	21755	
5	2011-12-09	2011-10-14	571255	22173	
6	2011-12-09	2011-10-17	C571499	23055	
7	2011-12-09	2011-03-18	547005	23177	
8	2011-12-09	2011-01-07	540480	22450	

```
-- 유저 별로 가장 큰 InvoiceDay를 찾아서 가장 최근 구매일로 저장
```

```
SELECT
```

```
CustomerID,
```

```
MAX(DATE(InvoiceDate)) AS InvoiceDay
```

```
FROM `modulabs-project-465302.modulabs_project.data`
WHERE CustomerID IS NOT NULL
GROUP BY CustomerID
```

행	CustomerID ▼	InvoiceDay ▼
1	12346	2011-01-18
2	12347	2011-12-07
3	12348	2011-09-25
4	12349	2011-11-21
5	12350	2011-02-02
6	12352	2011-11-03
7	12353	2011-05-19
8	12354	2011-04-21

```
-- 가장 최근 일자(most_recent_date)와 유저별 마지막 구매일(InvoiceDay)간의 차이를 구
SELECT
  CustomerID,
  EXTRACT(DAY FROM MAX(InvoiceDay) OVER () - InvoiceDay) AS recency
FROM (
  SELECT
    CustomerID,
    MAX(DATE(InvoiceDate)) AS InvoiceDay
  FROM `modulabs-project-465302.modulabs_project.data`
  GROUP BY CustomerID
);
```


행	CustomerID	recency
1	12407	49
2	12489	336
3	12577	35
4	12578	21
5	12581	39
6	12684	7
7	12712	22
8	12715	106

-- 지금까지의 결과를 user_r이라는 이름의 테이블로 저장

```
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.user_
```

```
WITH user_last_purchase AS (
```

```
  SELECT
```

```
    CustomerID,
```

```
    MAX(DATE(InvoiceDate)) AS InvoiceDay
```

```
  FROM `modulabs-project-465302.modulabs_project.data`
```

```
  WHERE CustomerID IS NOT NULL
```

```
  GROUP BY CustomerID
```

```
),
```

```
global_last_purchase AS (
```

```
  SELECT MAX(DATE(InvoiceDate)) AS most_recent_date
```

```
  FROM `modulabs-project-465302.modulabs_project.data`
```

```
)
```

```
SELECT
```

```
  u.CustomerID,
```

```
  u.InvoiceDay,
```

```
  g.most_recent_date,
```

```
  DATE_DIFF(g.most_recent_date, u.InvoiceDay, DAY) AS recency
```

```
FROM user_last_purchase u
```

```
CROSS JOIN global_last_purchase g
```

쿼리 결과

작업 정보 **결과** 실행 세부정보 실행 그래프

i 이 문으로 이름이 user_r인 새 테이블이 생성되었습니다.

저장소 프리뷰

```
-- Frequency
-- 1. 전체 거래 건수 계산
SELECT
  CustomerID,
  COUNT(DISTINCT InvoiceNo) AS purchase_cnt
FROM `modulabs-project-465302.modulabs_project.data`
WHERE CustomerID IS NOT NULL
GROUP BY CustomerID
```

행	CustomerID ▼	purchase_cnt ▼
1	12346	2
2	12347	7
3	12348	4
4	12349	1
5	12350	1
6	12352	8
7	12353	1
8	12354	1

```
-- 2. 구매한 아이템의 총 수량 계산
-- 각 고객 별로 구매한 아이템의 총 수량을 더해줌
SELECT
  CustomerID,
  SUM(Quantity) AS item_cnt
```

```
FROM `modulabs-project-465302.modulabs_project.data`
WHERE CustomerID IS NOT NULL
GROUP BY CustomerID
```

행	CustomerID	item_cnt
1	12346	0
2	12347	2446
3	12348	2332
4	12349	618
5	12350	196
6	12352	463
7	12353	20
8	12354	528

```
-- '1. 전체 거래 건수 계산'과 '2. 구매한 아이템의 총 수량 계산'의 결과를 합쳐서 user_rf라는
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.user_
```

```
WITH purchase_count AS (
  SELECT
    CustomerID,
    COUNT(DISTINCT InvoiceNo) AS purchase_cnt
  FROM `modulabs-project-465302.modulabs_project.data`
  WHERE CustomerID IS NOT NULL
  GROUP BY CustomerID
),
```

```
item_quantity AS (
  SELECT
    CustomerID,
    SUM(Quantity) AS item_cnt
  FROM `modulabs-project-465302.modulabs_project.data`
  WHERE CustomerID IS NOT NULL
```

```

GROUP BY CustomerID
)

SELECT
  pc.CustomerID,
  pc.purchase_cnt,
  iq.item_cnt
FROM purchase_count pc
JOIN item_quantity iq
ON pc.CustomerID = iq.CustomerID

```

The screenshot shows a database interface with a left sidebar containing a tree view of the database structure. The tree view includes 'dataset_test', 'modulabs_project', 'data', 'user_r', and 'user_rf'. The 'user_rf' table is highlighted in yellow. Below the tree view, there are buttons for '저장' (Save) and '프리뷰' (Preview). To the right of the tree view, there is a section titled '쿼리 결과' (Query Result) with tabs for '작업 정보' (Job Info), '결과' (Result), '실행 세부정보' (Execution Details), and '실행 그래프' (Execution Graph). The '결과' tab is selected, and it displays a message: '이 문으로 이름이 user_rf인 새 테이블이 생성되었습니다.' (A new table named user_rf is created by this statement.).

```

-- Monetary
-- 1. 고객별 총 지출액 계산
SELECT
  CustomerID,
  ROUND(SUM(Quantity * UnitPrice), 1) AS user_total
FROM `modulabs-project-465302.modulabs_project.data`
WHERE CustomerID IS NOT NULL
GROUP BY CustomerID

```

행	CustomerID	user_total
1	12346	0.0
2	12347	4302.2
3	12348	1437.2
4	12349	1398.1
5	12350	294.4
6	12352	1265.4
7	12353	89.0
8	12354	1045.5

-- 2. 고객별 평균 거래 금액 계산

-- 고객별 평균 거래 금액을 구하기 위해 1) data 테이블을 user_rf 테이블과 조인(LEFT JOIN)
 CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.user_`
 WITH user_base AS (

SELECT

r.CustomerID,

r.recency,

rf.purchase_cnt,

rf.item_cnt

FROM `modulabs-project-465302.modulabs_project.user_r` r

JOIN `modulabs-project-465302.modulabs_project.user_rf` rf

ON r.CustomerID = rf.CustomerID

),

user_total AS (

SELECT

CustomerID,

ROUND(SUM(Quantity * UnitPrice), 1) AS user_total

FROM `modulabs-project-465302.modulabs_project.data`

WHERE CustomerID IS NOT NULL

GROUP BY CustomerID

)

SELECT

ub.CustomerID,

```

ub.purchase_cnt,
ub.item_cnt,
ub.recency,
ut.user_total,
ROUND(ut.user_total / ub.purchase_cnt, 1) AS user_average
FROM user_base ub
LEFT JOIN user_total ut
ON ub.CustomerID = ut.CustomerID

```

```

-- RFM 통합 테이블 출력하기
SELECT *
FROM `modulabs-project-465302.modulabs_project.user_rfm`

```

행	CustomerID	purchase_cnt	item_cnt	recency	user_total	user_average
1	12713	1	493	0	779.5	779.5
2	14569	1	79	1	227.4	227.4
3	13436	1	76	1	196.9	196.9
4	13298	1	96	1	360.0	360.0
5	15520	1	314	1	343.5	343.5
6	15471	1	250	2	443.7	443.7
7	15195	1	1404	2	3861.0	3861.0
8	14204	1	72	2	150.6	150.6

```

-- RFM 고유한 유저의 수
SELECT COUNT(*) AS row_count

```

```
FROM `modulabs-project-465302.modulabs_project.user_rfm`
```

행	row_count ▼
1	4361

```
-- 추가 Feature 추출
```

```
-- 구매하는 제품의 다양성
```

```
-- 1) 고객 별로 구매한 상품들의 고유한 수를 계산합니다. 높은 숫자가 나오는 것은 해당 고객
```

```
-- 이후 2) user_rfm 테이블과 결과를 합치고, 이를 3) user_data라는 이름의 테이블에 저장
```

```
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.user_`  
WITH unique_products AS (
```

```
  SELECT
```

```
    CustomerID,
```

```
    COUNT(DISTINCT StockCode) AS unique_products
```

```
  FROM `modulabs-project-465302.modulabs_project.data`
```

```
  WHERE CustomerID IS NOT NULL
```

```
  GROUP BY CustomerID
```

```
)
```

```
SELECT
```

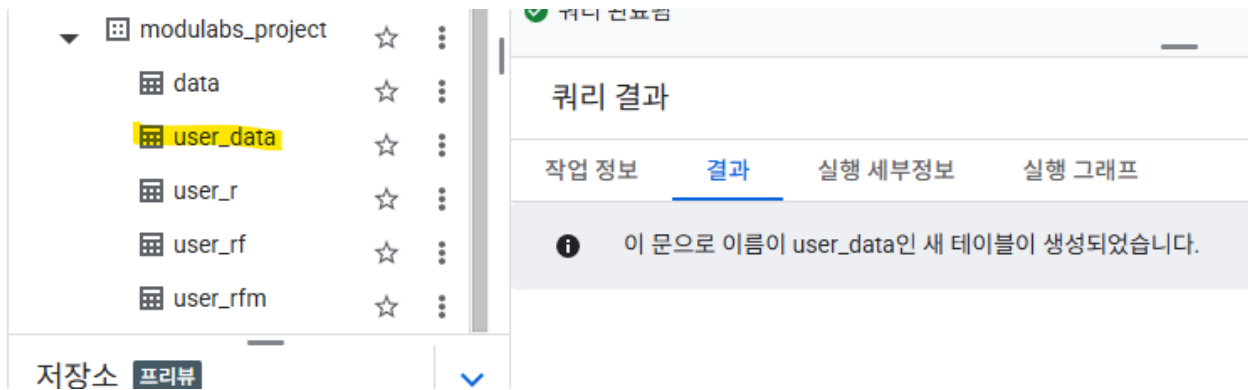
```
  ur.*,
```

```
  up.unique_products
```

```
FROM `modulabs-project-465302.modulabs_project.user_rfm` AS ur
```

```
JOIN unique_products AS up
```

```
ON ur.CustomerID = up.CustomerID
```



```
-- 평균 구매 주기
-- 평균 구매 소요 일수를 계산하고, 그 결과를 user_data에 통합
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.user_
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
  FROM (
    -- (1) 구매와 구매 사이에 소요된 일수
    SELECT
      CustomerID,
      DATE_DIFF(DATE(InvoiceDate),
        LAG(DATE(InvoiceDate)) OVER (PARTITION BY CustomerID ORDER BY
          DAY) AS interval_
      FROM `modulabs-project-465302.modulabs_project.data`
      WHERE CustomerID IS NOT NULL
    )
    GROUP BY CustomerID
  )

  SELECT
    u.*,
```



```

pi.average_interval
FROM `modulabs-project-465302.modulabs_project.user_data` AS u
LEFT JOIN purchase_intervals AS pi
ON u.CustomerID = pi.CustomerID

```

 이 문으로 이름이 user_data인 테이블이 교체되었습니다.

```

-- 구매 취소 경향성
-- 취소 빈도와 취소 비율을 계산하고 그 결과를 user_data에 통합
CREATE OR REPLACE TABLE `modulabs-project-465302.modulabs_project.user_

WITH user_rfm AS (
  SELECT *
  FROM `modulabs-project-465302.modulabs_project.user_rfm`
),
TransactionInfo AS (
  SELECT
    CustomerID,
    COUNT(DISTINCT InvoiceNo) AS total_transactions,
    COUNT(DISTINCT CASE WHEN InvoiceNo LIKE 'C%' THEN InvoiceNo END) AS
  FROM `modulabs-project-465302.modulabs_project.data`
  WHERE CustomerID IS NOT NULL
  GROUP BY CustomerID
)

SELECT
  u.*,
  t.total_transactions,
  t.cancel_frequency,
  ROUND(t.cancel_frequency / t.total_transactions * 100, 2) AS cancel_rate
FROM user_rfm u
LEFT JOIN TransactionInfo t
ON u.CustomerID = t.CustomerID

```

i 이 문으로 이름이 user_data인 테이블이 교체되었습니다.

-- 최종적으로 user_data를 출력

SELECT *

FROM `modulabs-project-465302.modulabs_project.user_data`

행	CustomerID	purchase_cnt	item_cnt	recency	user_total	user_average
1	14703	1	183	14	318.2	31
2	17385	1	197	14	256.1	25
3	12890	1	252	24	372.9	37
4	12552	1	85	39	317.8	31
5	12588	1	52	39	174.9	17
6	16127	1	281	39	606.0	60
7	18067	1	189	40	396.3	39
8	14585	1	91	50	157.1	15

• 파이썬으로 열어보기

```
[1]: import pandas as pd
user_data = pd.read_csv('./data/user_data.csv')

[2]: user_data.head()
```

	CustomerID	purchase_cnt	item_cnt	recency	user_total	user_average	total_transactions	cancel_frequency	cancel_rate
0	15857	1	308	18	297.0	297.0	1	0	0.0
1	12445	1	60	22	77.4	77.4	1	0	0.0
2	15556	1	236	24	279.8	279.8	1	0	0.0
3	16127	1	281	39	606.0	606.0	1	0	0.0
4	15096	1	120	42	219.4	219.4	1	0	0.0