

# 고객을 세그멘테이션하자 [프로젝트]

## 11-2. 데이터 불러오기

### 데이터 살펴보기

- 테이블에 있는 10개의 행만 출력하기

```
select * from white-might-456101-t9.modulabs_project.data limit 10;
```

쿼리 결과										결과 저장	다음에서 열기	
작업 정보	결과	차트	JSON	실행 세부정보	실행 그래프							
행	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country				
1	536363	89123A	WHITE HANGING HEART TAIL...	6	2010-12-01 09:26:00 UTC	2.05	17850	United Kingdom				
2	536363	71053	WHITE METAL LANTERN	6	2010-12-01 09:26:00 UTC	3.39	17850	United Kingdom				
3	536363	84406B	CREAM CUPID HEARTS COAT	6	2010-12-01 09:26:00 UTC	2.75	17850	United Kingdom				
4	536363	84209G	KNITTED UNION FLAG HOT WIL...	6	2010-12-01 09:26:00 UTC	3.39	17850	United Kingdom				
5	536363	84207E	RED WOOLLY HOTTIE WHITE H...	6	2010-12-01 09:26:00 UTC	3.39	17850	United Kingdom				
6	536363	22752	SET 7 BABUSHKA NESTING BO...	2	2010-12-01 09:26:00 UTC	7.65	17850	United Kingdom				
7	536363	21750	GLASS STAR FROSTED TULIP...	6	2010-12-01 09:26:00 UTC	4.25	17850	United Kingdom				
8	536363	22833	HAND WARMER UGGY JACK	6	2010-12-01 09:26:00 UTC	1.85	17850	United Kingdom				
9	536363	22832	HAND WARMER RED PINKA D...	6	2010-12-01 09:26:00 UTC	1.85	17850	United Kingdom				
10	536367	84479	ASSORTED COLOUR BIRD ORN...	32	2010-12-01 09:34:00 UTC	1.49	13047	United Kingdom				

- 전체 데이터는 몇 행으로 구성되어 있는지 확인하기

```
select count(*) as total_count from white-might-456101-t9.modulabs_project.data;
```

쿼리 결과

결과 저장

다음에서 열기

작업 정보

결과

차트

JSON

실행 세부정보

실행 그래프

행	total_count	
1	541909	

### 데이터 수 세기

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

```
SELECT
COUNT(InvoiceNo) AS InvoiceNo_count,
COUNT(StockCode) AS StockCode_count,
COUNT(Description) AS Description_count,
COUNT(Quantity) AS Quantity_count,
COUNT(InvoiceDate) AS InvoiceDate_count,
COUNT(UnitPrice) AS UnitPrice_count,
COUNT(CustomerID) AS CustomerID_count,
COUNT(Country) AS Country_count
FROM
`white-might-456101-t9.modulabs_project.data`;
```

쿼리 결과

작업 정보

결과

차트

JSON

실행 세부정보

실행 그래프

행	InvoiceNo_count	StockCode_count	Description_count	Quantity_count	InvoiceDate_count	UnitPrice_count	CustomerID_count	Country_count
1	541909	541909	540455	541909	541909	541909	406829	541909

## 11-4. 데이터 전처리 방법(1): 결측치 제거

## 컬럼 별 누락된 값의 비율 계산

- 각 컬럼 별 누락된 값의 비율을 계산
  - 각 컬럼에 대해서 누락 값을 계산한 후, 계산된 누락 값을 **UNION ALL**을 통해 합치기

```
SELECT column_name, ROUND((total - column_value) / total * 100, 2)
FROM
(
  SELECT 'InvoiceNo' AS column_name, COUNT(InvoiceNo) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
  SELECT 'StockCode' AS column_name, COUNT(StockCode) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
  SELECT 'Description' AS column_name, COUNT(Description) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
  SELECT 'Quantity' AS column_name, COUNT(Quantity) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
  SELECT 'InvoiceDate' AS column_name, COUNT(InvoiceDate) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
  SELECT 'UnitPrice' AS column_name, COUNT(UnitPrice) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
  SELECT 'CustomerID' AS column_name, COUNT(CustomerID) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
  SELECT 'Country' AS column_name, COUNT(Country) AS column_value, COUNT(*) AS total FROM white-might-456101-t9.modulabs_project.data
) AS column_data;
```

쿼리 결과

작업 정보	결과	차트	JSON	실행 세부정보
행	column_name	fd_		
1	InvoiceDate	0.0		
2	CustomerID	24.93		
3	Quantity	0.0		
4	InvoiceNo	0.0		
5	StockCode	0.0		
6	UnitPrice	0.0		
7	Description	0.27		
8	Country	0.0		

## 결측치 처리 전략

- **StockCode = '85123A'** 의 **Description** 을 추출하는 쿼리문을 작성하기

```
select distinct Description from white-might-456101-t9.modulabs_project.data where StockCode = '85123A';
```

쿼리 결과

작업 정보	결과	차트	JSON
행	Description		
1	WHITE HANGING HEART T-LUG...		
2	?		
3	wrongly marked carton 22804		
4	CREAM HANGING HEART T-LUG...		

## 결측치 처리

- **DELETE** 구문을 사용하며, **WHERE** 절을 통해 데이터를 제거할 조건을 제시

```
delete FROM `white-might-456101-t9.modulabs_project.data`
WHERE CustomerID IS NULL OR Description IS NULL;
```

쿼리 결과			
작업 정보	결과	실행 세부정보	실행 그래프
<div> 이 문으로 data의 행 135,080개가 삭제되었습니다. </div>			

## 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 double_count
FROM (
  SELECT *
  FROM `white-might-456101-t9.modulabs_project.data`
  GROUP BY InvoiceNo, StockCode, Description, Quantity, InvoiceDate, UnitPrice, CustomerID, Country
  HAVING COUNT(*) > 1
);
```

쿼리 결과		
작업 정보	결과	실행 세부정보
행	double_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 `white-might-456101-t9.modulabs_project.data` AS
SELECT DISTINCT *
FROM `white-might-456101-t9.modulabs_project.data`;
```

쿼리 결과			
작업 정보	결과	실행 세부정보	실행 그래프
<div> 이 문으로 이름이 data인 테이블이 교체되었습니다. </div>			

## 11-6. 데이터 전처리(3): 오류값 처리

### InvoiceNo 살펴보기

- 고유(unique)한 InvoiceNo의 개수를 출력하기

```
select count(distinct InvoiceNo)
from white-might-456101-t9.modulabs_project.data;
```

쿼리 결과

작업 정보	결과	차트
행	fo_	
1	22190	

- 고유한 **InvoiceNo** 를 앞에서부터 100개를 출력하기

```
select distinct InvoiceNo
from white-might-456101-t9.modulabs_project.data
limit 100;
```

쿼리 결과

작업 정보	결과	차트	JSON	실행 세부정보	실행 그래프
행	InvoiceNo				
1	541431				
2	C541433				
3	537626				
4	542237				
5	549222				
6	556201				
7	562032				
8	573511				
9	581180				
10	539318				
11	541998				
12	548955				
13	568172				
14	577609				

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

```
select *
from white-might-456101-t9.modulabs_project.data
where InvoiceNo LIKE 'C%'
limit 100;
```

쿼리 결과

작업 정보	결과	차트	JSON	실행 세부정보	실행 그래프			
행	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
30	C565050	2243	SET OF TEA COFFEE SUGAR TL...	-2	2011-08-31 17:10:00 UTC	4.95	12381	Norway
31	C565050	851098	WHITE TEA COFFEE SUGAR JA...	-1	2011-08-31 17:10:00 UTC	1.95	12381	Norway
32	C565050	21791	VINTAGE HEADS AND TAILS C...	-12	2011-08-31 17:10:00 UTC	1.25	12381	Norway
33	C565050	23198	PANTRY MAGNETIC SHOPPIN...	-10	2011-08-31 17:10:00 UTC	1.45	12381	Norway
34	C541586	21218	RED SPOTTY BISCUIT TIN	-3	2011-01-19 14:36:00 UTC	3.75	12383	Belgium
35	C567540	23295	SET OF 12 MINI LOAF BAKING ...	-2	2011-09-21 10:00:00 UTC	0.83	12384	Switzerland
36	C567540	23173	REGENCY TEAPOT ROSES	-1	2011-09-21 10:00:00 UTC	9.95	12384	Switzerland
37	C567540	84279P	CHERRY BLOSSOM DECORATI...	-2	2011-09-21 10:00:00 UTC	3.75	12384	Switzerland
38	C540152	22666	RECIFE BOX PANTRY YELLOW ...	-1	2011-01-05 11:21:00 UTC	2.95	12395	Belgium
39	C544419	85232A	SET/3 POLKADOT STACKING T...	-3	2011-02-18 15:16:00 UTC	4.95	12395	Belgium
40	C579432	21599	STRAWBERRY LUNCH BOX WL...	-1	2011-11-24 12:20:00 UTC	2.55	12395	Belgium
41	C572187	P057	POSTAGE	-1	2011-10-01 11:00:00 UTC	18.0	12403	Denmark
42	C562728	22849	BREAD BIN DINER STYLE MINT	-1	2011-08-09 09:41:00 UTC	14.95	12406	Denmark
43	C562728	22846	BREAD BIN DINER STYLE RED	-4	2011-08-09 09:41:00 UTC	14.95	12406	Denmark
44	C562728	22848	BREAD BIN DINER STYLE FRANK	-4	2011-08-09 09:41:00 UTC	14.95	12406	Denmark

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

```
SELECT ROUND(SUM(CASE WHEN InvoiceNo like 'C%' THEN 1 ELSE 0 END)/ count(InvoiceNo) * 100, 1) as Canceled_Percent
from white-might-456101-t9.modulabs_project.data
```

쿼리 결과	
작업 정보	결과
행	Canceled_Percent
1	2.2

## StockCode 살펴보기

- 고유한 StockCode 의 개수를 출력하기

```
select count(distinct StockCode)
from white-might-456101-t9.modulabs_project.data;
```

쿼리 결과	
작업 정보	결과
행	fo_
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 white-might-456101-t9.modulabs_project.data
group by StockCode
ORDER BY sell_cnt DESC
limit 10;
```

쿼리 결과	
작업 정보	결과
행	StockCode
1	85123A
2	22423
3	85099B
4	47566
5	84879
6	20725
7	22720
8	POST
9	22197
10	23203

- StockCode 의 컬럼에 있던 값 중에서 숫자를 제외한 문자만 남기고 문자가 몇 자리 수 인지 세고
  - 숫자가 0~1개인 값들에는 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 white-might-456101-t9.modulabs_project.data
)
WHERE number_count between 0 and 1;
```

쿼리 결과

작업 정보	결과	자트	JSON	실행 세부정보	실행 그래프
행	StockCode	number_count			
1	POST	0			
2	M	0			
3	CZ	1			
4	D	0			
5	BANK CHARGES	0			
6	PADS	0			
7	DOT	0			
8	CRUK	0			

- **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)

```
-- 숫자가 0~1개인 값
WITH no_char_codes AS (
  SELECT DISTINCT StockCode
  FROM (
    SELECT StockCode,
      LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS number_count
    FROM `white-might-456101-t9.modulabs_project.data`
  )
  WHERE number_count BETWEEN 0 AND 1
),

-- no_char_codes에서 distinct StockCode에 출력된 결과만을 포함하는 전체 데이터수
filter_no_char_codes AS (
  select count(*) as row_count
  from `white-might-456101-t9.modulabs_project.data`
  where StockCode in (select StockCode from no_char_codes)
),

-- data 테이블의 전체 데이터 행 수
total AS (
  select count(*) as total_count from `white-might-456101-t9.modulabs_project.data`
)

select ROUND( row_count / total.total_count * 100, 2) as no_char_stockcode_percent
from filter_no_char_codes, total;
```

쿼리 결과

작업 정보	결과	자트	JSON	실행 세부정보	실행 그래프
행	no_char_stockcode_percent				
1	0.48				

- 제품과 관련되지 않은 거래 기록을 제거하기

```
delete from `white-might-456101-t9.modulabs_project.data`
where StockCode in (
  SELECT DISTINCT StockCode
  FROM (
    SELECT StockCode,
      LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS number_count
    FROM `white-might-456101-t9.modulabs_project.data`
  )
  WHERE number_count > 1
);
```

```

)
WHERE number_count BETWEEN 0 AND 1
);

```

#### 쿼리 결과

작업 정보   **결과**   실행 세부정보   실행 그래프

**i** 이 문으로 data의 행 1,915개가 삭제되었습니다.

## Description 살펴보기

- 고유한 Description 별 출현 빈도를 계산하고 상위 30개를 출력하기

```

SELECT DISTINCT Description,COUNT(*) AS description_cnt
FROM `white-might-456101-t9.modulabs_project.data`
WHERE REGEXP_CONTAINS(Description, r'[a-z]')
group by Description
order by description_cnt DESC
limit 30;

```

#### 쿼리 결과

작업 정보	결과	자트	JSON	실행 세부정보	실행 그래프
행	Description	description_cnt			
1	BAG 125g SWIRLY MARBLES	252			
2	3 TRADITIONAL BISCUIT OUTT...	207			
3	BAG 250g SWIRLY MARBLES	201			
4	POLYESTER FILLER PAD 40x40...	185			
5	POLYESTER FILLER PAD 45x45...	135			
6	BAG 500g SWIRLY MARBLES	115			
7	Next Day Carriage	80			
8	FRENCH BLUE METAL DOOR SL...	75			
9	POLYESTER FILLER PAD 45x30...	37			
10	POLYESTER FILLER PAD 30CM...	26			
11	ESSENTIAL BALM 3.5g TIN IN ...	18			
12	NUMBER TILE COTTAGE GARD...	12			
13	FOLK ART GREETING CARD,pa...	10			
14	THE KING GIFT BAG 25x24x12...	7			
15	NUMBER TILE VINTAGE FONT ...	7			
16	POLYESTER FILLER PAD 65CM...	5			
17	High Resolution Image	3			
18	FLOWERS HANDBAG blue and ...	3			
19	POLYESTER FILLER PAD 60x40...	1			

- 서비스 관련 정보를 포함하는 행들을 제거하기

```

delete from `white-might-456101-t9.modulabs_project.data`
where Description = 'Next Day Carriage' or Description = 'High Resolution Image';

```

#### 정답코드

```

DELETE
FROM project_name.modulabs_project.data
WHERE
Description LIKE '%Next Day Carriage%'
OR Description LIKE '%High Resolution Image%';

```

#### 쿼리 결과

작업 정보 **결과** 실행 세부정보 실행 그래프

이 문으로 data의 행 83개가 삭제되었습니다.

지정된 문자열이 아니고 "서비스 관련 정보를 포함하는" 이므로 %연산자 사용이 맞다고 봄.

- 대소문자를 혼합하고 있는 데이터를 대문자로 표준화 하기

```
CREATE OR REPLACE TABLE white-might-456101-t9.modulabs_project.data AS
SELECT
  * EXCEPT (Description),
  UPPER(Description) AS Description
FROM white-might-456101-t9.modulabs_project.data
```

#### 쿼리 결과

작업 정보 **결과** 실행 세부정보 실행 그래프

이 문으로 이름이 data인 테이블이 교체되었습니다.

## UnitPrice 살펴보기

- UnitPrice의 최솟값, 최댓값, 평균을 구하기

```
select min(UnitPrice) as min_price, max(UnitPrice) as max_price, avg(UnitPrice) as avg_price
from `white-might-456101-t9.modulabs_project.data`
```

#### 쿼리 결과

작업 정보	결과	자트	JSON	실행 세부정보	실행 그래프
행	min_price	max_price	avg_price		
1	0.0	649.5	2.904956757406...		

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

```
SELECT count(Quantity) AS cnt_quantity, min(Quantity) AS min_quantity, max(Quantity) AS max_quantity, avg(Quantity) AS avg_quar
FROM `white-might-456101-t9.modulabs_project.data`
WHERE UnitPrice = 0
```

#### 쿼리 결과

작업 정보	결과	자트	JSON	실행 세부정보	실행 그래프
행	cnt_quantity	min_quantity	max_quantity	avg_quantity	
1	33	1	12540	420.5151515151...	

- UnitPrice = 0 를 제거하고 일관된 데이터셋을 유지하기





쿼리 결과		
작업 정보 <b>결과</b> 차트   JSON   실행		
행	most_recent_date	
1	2011-12-09	

정답코드

```
SELECT
  MAX(DATE(InvoiceDate)) OVER () AS most_recent_date,
  DATE(InvoiceDate) AS InvoiceDay,
  *
FROM project_name.modulabs_project.data;
```

왜 over구문인지? → 꼭 필요한 구문은 아니나 모든 행에 추가해서 보여줄 경우 사용.

📌 차이를 한 줄로 말하면?

- **MAX()** → 전체 요약
- **MAX() OVER ()** → 전체 요약값을 모든 행에 추가해서 보여줌

OVER() 구문을 쓰지 않았다고 해서 오답 코드는 아니라고 봄.

- 유저 별로 가장 큰 InvoiceDay를 찾아서 가장 최근 구매일로 저장하기

```
SELECT
  CustomerID, MAX(DATE(InvoiceDate)) AS InvoiceDay
FROM `white-might-456101-t9.modulabs_project.data`
group by CustomerID
```

쿼리 결과		
작업 정보 <b>결과</b> 차트   JSON   실행		
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
9	12355	2011-05-09

- 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
FROM (
  SELECT
    CustomerID,
    MAX(DATE(InvoiceDate)) AS InvoiceDay
  FROM `white-might-456101-t9.modulabs_project.data`
  GROUP BY CustomerID
);
```

← 쿼리 결과

작업 정보	결과	차트	JSON	실행 세부정보	실행 그래프
행	CustomerID	recency			
1	12399	119			
2	12448	44			
3	12700	7			
4	12718	71			
5	13369	354			
6	13436	1			
7	13606	29			
8	13823	18			
9	14142	373			
10	14196	106			

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

```
CREATE OR REPLACE TABLE `white-might-456101-t9.modulabs_project.user_r` AS
WITH rd AS (
  SELECT
    CustomerID,
    MAX(DATE(InvoiceDate)) AS InvoiceDay
  FROM `white-might-456101-t9.modulabs_project.data`
  GROUP BY CustomerID
),
rd2 AS (
  SELECT
    MAX(DATE(InvoiceDate)) AS most_recent_date
  FROM `white-might-456101-t9.modulabs_project.data`
)
SELECT
  rd.CustomerID,
  rd.InvoiceDay,
  rd2.most_recent_date,
  DATE_DIFF(rd2.most_recent_date, rd.InvoiceDay, DAY) AS recency
FROM rd, rd2;
```

쿼리 결과

작업 정보	결과	실행 세부정보	실행 그래프
<p><b>i</b> 이 문으로 이름이 user_r인 테이블이 교체되었습니다.</p>			

## Frequency

- 고객마다 고유한 InvoiceNo의 수를 세어보기

```
SELECT
  CustomerID,
  count(distinct invoiceNo) AS purchase_cnt
FROM `white-might-456101-t9.modulabs_project.data`
group by CustomerID;
```

쿼리 결과			
작업 정보		결과	차트
명	CustomerID	purchase_cnt	
1	12457	1	
2	12847	1	
3	13985	1	
4	15107	1	
5	16818	1	
6	12444	1	
7	13014	1	
8	14410	1	
9	15904	1	
10	15902	1	

- 각 고객 별로 구매한 아이템의 총 수량 더하기

```
SELECT
  CustomerID, SUM(Quantity) AS item_cnt
FROM `white-might-456101-t9.modulabs_project.data`
group by CustomerID;
```

쿼리 결과			
작업 정보		결과	차트
명	CustomerID	item_cnt	
1	12346	0	
2	12547	2458	
3	12540	2332	
4	12349	630	
5	12350	196	
6	12352	463	
7	12353	20	
8	12354	530	
9	12355	240	
10	13356	1873	

쿼리 결과			
작업 정보		결과	차트
명	CustomerID	item_cnt	
16	12507	1	
17	12431	1	
18	14110	1	
19	17560	1	
20	13081	4	
21	13239	1	
22	12415	2	
23	14646	4	
24	13256	1	

- 전체 거래 건수와 구매한 아이템의 총 수량 계산의 결과를 합쳐서 **user\_rf** 라는 이름의 테이블에 저장하기

```
CREATE OR REPLACE TABLE white-might-456101-t9.modulabs_project.user_rf AS
WITH purchase_cnt AS (
  SELECT
    CustomerID,
    count(distinct invoiceNo) AS purchase_cnt
  FROM `white-might-456101-t9.modulabs_project.data`
  group by CustomerID
),

-- (2) 구매한 아이템 총 수량 계산
item_cnt AS (
  SELECT
    CustomerID, SUM(Quantity) AS item_cnt
  FROM `white-might-456101-t9.modulabs_project.data`
  group by CustomerID
)
```

```
-- 기존의 user_r에 (1)과 (2)를 통합
SELECT
  pc.CustomerID,
  pc.purchase_cnt,
  ic.item_cnt,
  ur.rececy
FROM purchase_cnt AS pc
JOIN item_cnt AS ic
  ON pc.CustomerID = ic.CustomerID
JOIN `white-might-456101-t9.modulabs_project.user_r` AS ur
  ON pc.CustomerID = ur.CustomerID;
```

쿼리 결과 결과 저장 다

작업 정보 **결과** 실행 세부정보 실행 그래프

**i** 이 문으로 이름이 user\_rf인 새 테이블이 생성되었습니다.

## Monetary

- 고객별 총 지출액 계산 (소수점 첫째 자리에서 반올림)

```
SELECT CustomerID,
  ROUND(SUM(UnitPrice), 1) AS user_total
FROM `white-might-456101-t9.modulabs_project.data`
group by CustomerID;
```

쿼리 결과			
작업 정보	결과	차트	JSON
실행 세부정보	실행 그래프		
행	CustomerID	user_total	
1	12346	2.1	
2	12347	481.2	
3	12348	18.7	
4	12349	305.1	
5	12350	25.3	
6	12352	330.5	
7	12353	24.3	
8	12354	261.2	
9	12355	54.7	

- 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 white-might-456101-t9.modulabs_project.user_rfm AS

SELECT
  rf.CustomerID AS CustomerID,
  rf.purchase_cnt,
  rf.item_cnt,
```

```

rf.recency,
ut.user_total,
round(ut.user_total / rf.purchase_cnt,2) AS user_average
FROM white-might-456101-t9.modulabs_project.user_rf rf
LEFT JOIN (
SELECT CustomerID,
ROUND(SUM(UnitPrice), 1) AS user_total
FROM `white-might-456101-t9.modulabs_project.data`
group by CustomerID
) ut
ON rf.CustomerID = ut.CustomerID;

```

쿼리 결과	
작업 정보	결과
실행 세부정보	
실행 그래프	
이 문으로 이름이 user_rfm인 테이블이 교체되었습니다.	

## RFM 통합 테이블 출력하기

- 최종 user\_rfm 테이블을 출력하기

```
select * FROM `white-might-456101-t9.modulabs_project.user_rfm`
```

쿼리 결과							
작업 정보	결과	자트	JSON	실행 세부정보	실행 그래프		
행	CustomerID	purchase_cnt	item_cnt	recency	user_total	user_average	
1	12713	1	505	0	77.1	77.1	
2	13298	1	96	1	7.5	7.5	
3	15520	1	314	1	31.0	31.0	
4	13436	1	76	1	70.0	70.0	
5	14569	1	79	1	47.0	47.0	
6	15471	1	256	2	146.4	146.4	
7	15195	1	1404	2	2.8	2.8	
8	14204	1	72	2	99.2	99.2	
9	15318	1	642	3	19.3	19.3	

고유한 CustomerID 수

쿼리 결과	
작업 정보	결과
자트	JSON
실행 세부정보	실행 그래프
행	f0_
1	4362

## 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 `white-might-456101-t9.modulabs_project.user_data` AS
WITH unique_products AS (
  SELECT
    CustomerID,
    COUNT(DISTINCT StockCode) AS unique_products
  FROM `white-might-456101-t9.modulabs_project.data`
  GROUP BY CustomerID
)
SELECT ur.*, up.* EXCEPT (CustomerID)
FROM white-might-456101-t9.modulabs_project.user_rfm AS ur
JOIN unique_products AS up
ON ur.CustomerID = up.CustomerID;
```

쿼리 결과	
작업 정보	결과
	실행 세부정보
	실행 그래프
<div> <div></div> <div>이 문으로 이름이 <code>user_data</code>인 새 테이블이 생성되었습니다.</div> </div>	

## 2. 평균 구매 주기

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

```
CREATE OR REPLACE TABLE `white-might-456101-t9.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
  FROM (
    -- (1) 구매와 구매 사이에 소요된 일수
    SELECT
      CustomerID,
      DATE_DIFF(InvoiceDate, LAG(InvoiceDate) OVER (PARTITION BY CustomerID ORDER BY InvoiceDate), DAY) AS interval_
    FROM
      `white-might-456101-t9.modulabs_project.data`
    WHERE CustomerID IS NOT NULL
  )
  GROUP BY CustomerID
)
SELECT u.*, pi.* EXCEPT (CustomerID)
FROM `white-might-456101-t9.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
<div> <span>❗</span> 이 문으로 이름이 user_data인 테이블이 교체되었습니다. </div>			

### 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 `white-might-456101-t9.modulabs_project.user_data` AS
WITH TransactionInfo AS (
  SELECT
    CustomerID,
    count(case when InvoiceNo LIKE 'C%' then 1 end) AS cancel_frequency,
    COUNT(InvoiceNo) AS total_transactions
  FROM `white-might-456101-t9.modulabs_project.data`
  group by CustomerID
)

SELECT u.*, t.* EXCEPT(CustomerID), round( (t.cancel_frequency / t.total_transactions)*100, 2) AS cancel_rate
FROM `white-might-456101-t9.modulabs_project.user_data` AS u
LEFT JOIN TransactionInfo AS t
ON u.CustomerID = t.CustomerID;
```

쿼리 결과			
작업 정보	결과	실행 세부정보	실행 그래프
<div> <span>❗</span> 이 문으로 이름이 user_data인 테이블이 교체되었습니다. </div>			

#### 정답코드

```
CREATE OR REPLACE TABLE `project_name.modulabs_project.user_data` AS

WITH TransactionInfo AS (
  SELECT
    CustomerID,
    COUNT(DISTINCT InvoiceNo) AS total_transactions,
    COUNT(DISTINCT CASE WHEN InvoiceNo LIKE 'C%' THEN InvoiceNo END) AS cancel_frequency
  FROM `project_name.modulabs_project.data`
  GROUP BY CustomerID
)

SELECT u.*, t.* EXCEPT(CustomerID), ROUND(cancel_frequency/total_transactions, 2) AS cancel_rate
FROM `project_name.modulabs_project.user_data` AS u
```



