

강상현_고객을 세그먼테이션하자 [프로젝트]

11-2. 데이터 불러오기

데이터 살펴보기

- 테이블에 있는 10개의 행만 출력하기

```
SELECT *
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
LIMIT 10
```

[결과 이미지를 넣어주세요]

쿼리 결과										
작업 정보		결과		시각화		JSON		실행 세부정보		실행 그레프
행	InvoiceNo	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country		
1	536365	85123A	WHITE HANGING HEART TUS...	6	2010-12-01 08:26:00 UTC	2.59	17850	United Kingdom		
2	536365	71053	WHITE METAL LANTERN	6	2010-12-01 08:26:00 UTC	3.39	17850	United Kingdom		
3	536365	844068	CREAM CUPID HEARTS COAT H...	8	2010-12-01 08:26:00 UTC	2.75	17850	United Kingdom		
4	536365	840296	KNITTED UNION FLAG HOT WA...	6	2010-12-01 08:26:00 UTC	3.39	17850	United Kingdom		
5	536365	840296	RED WOOLLY HOTTIE WHITE H...	6	2010-12-01 08:26:00 UTC	3.39	17850	United Kingdom		
6	536365	22752	SET 7 BABUSHKA NESTING BO...	2	2010-12-01 08:26:00 UTC	7.65	17850	United Kingdom		
7	536365	21730	GLASS STAR FROSTED TLIGHT...	6	2010-12-01 08:26:00 UTC	4.25	17850	United Kingdom		
8	536366	22633	HAND WARMER UNION JACK	6	2010-12-01 08:26:00 UTC	1.85	17850	United Kingdom		
9	536366	22632	HAND WARMER RED POLKA DOT	6	2010-12-01 08:26:00 UTC	1.85	17850	United Kingdom		
10	536367	84079	ASSORTED COLOUR BIRD ORN...	32	2010-12-01 08:34:00 UTC	1.69	13047	United Kingdom		

- 전체 데이터는 몇 행으로 구성되어 있는지 확인하기

```
SELECT COUNT(*)
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보		결과
행	f0_	
1	541909	

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 `coral-bucksaw-482803-p2.modulabs_project.data`
```

[결과 이미지를 넣어주세요]

쿼리 결과									
작업 정보		결과		시각화		JSON		실행 세부정보	
행	COUNT_InvoiceNo	COUNT_StockCode	COUNT_Description	COUNT_Quantity	COUNT_InvoiceD...	COUNT_UnitPrice	COUNT_Custome...	COUNT_Country	
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 ALL을 통해 합치기

```
SELECT
    'InvoiceNo' AS column_name,
    ROUND(SUM(CASE WHEN InvoiceNo IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
UNION ALL
SELECT
    'StockCode' AS column_name,
    ROUND(SUM(CASE WHEN StockCode IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
UNION ALL
SELECT
    'Description' AS column_name,
    ROUND(SUM(CASE WHEN Description IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`

UNION ALL
SELECT
    'Quantity' AS column_name,
    ROUND(SUM(CASE WHEN Quantity IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`

UNION ALL
SELECT
    'InvoiceDate' AS column_name,
    ROUND(SUM(CASE WHEN InvoiceDate IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`

UNION ALL
SELECT
    'UnitPrice' AS column_name,
    ROUND(SUM(CASE WHEN UnitPrice IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`

UNION ALL
SELECT
    'CustomerID' AS column_name,
    ROUND(SUM(CASE WHEN CustomerID IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`

UNION ALL
SELECT
    'Country' AS column_name,
    ROUND(SUM(CASE WHEN Country IS NULL THEN 1 ELSE 0 END) / COUNT(*) * 100,2) AS missing_percentage
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
```

[결과 이미지를 넣어주세요]

쿼리 결과

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

결측치 처리 전략

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

```
SELECT DISTINCT Description
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
WHERE StockCode = '85123A'
```

[결과 이미지를 넣어주세요]

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

결측치 처리

- DELETE 구문을 사용하여, WHERE 절을 통해 데이터를 제거할 조건을 제시

```
DELETE
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
WHERE Description IS NULL
OR CustomerID IS NULL
```

[결과 이미지를 넣어주세요]

쿼리 결과

작업 정보	결과	실행 세부정보	실행 그래프
❶ 이 문으로 data의 행 1,454개가 삭제되었습니다.			

>Description IS NULL 제거 값

작업 정보	결과	실행 세부정보	실행 그래프
❶ 이 문으로 data의 행 133,626개가 삭제되었습니다.			

>CustomerID IS NULL 제거 값
최종 135,080개가 제거 됨

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 *
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY InvoiceNo, StockCode, Description, Quantity, InvoiceDate, UnitPrice, CustomerID, Country
HAVING COUNT(*) > 1
```

[결과 이미지를 넣어주세요]

쿼리 결과						
작업 정보		결과	시각화	JSON	실행 세부정보	실행 그래프
행	InvoiceNo	StockCode	Description	Quantity	Invoic	
1	571034	23494	VINTAGE DOILY DELUXE SEWIN...	3	2011-	
2	571034	23239	SET OF 4 KNICK KNACK TINS P...	6	2011-	
3	538826	22749	FELTCRAFT PRINCESS CHARLO...	1	2010-	
4	577228	22435	SET OF 9 HEART SHAPED BALL...	1	2011-	
5	577228	23156	SET OF 5 MINI GROCERY MAG...	1	2011-	

페이지당 결과 수: 50 1 – 50 (전체 4837행) |< < > >|

중복값 처리

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

```
CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.data` AS
SELECT DISTINCT *
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
```

[결과 이미지를 넣어주세요]

```
113 CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.data` AS
114 SELECT DISTINCT *
115 FROM `coral-bucksaw-482803-p2.modulabs_project.data`
```

116 쿼리 인료됨

주문형 처리 할당량 사용 중

쿼리 결과

작업 정보 결과 실행 세부정보 실행 그래프

❶ 이 문으로 이름이 data인 테이블이 교체되었습니다.

쿼리 결과		
작업 정보		결과
행	f0_	
1	401604	

11-6. 데이터 전처리(3): 오류값 처리

InvoiceNo 살펴보기

- 고유(unique)한 InvoiceNo의 개수를 출력하기

```
SELECT COUNT(DISTINCT InvoiceNo)
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보	결과	시각화
행	f0_	
1	22190	

- 고유한 InvoiceNo를 앞에서부터 100개를 출력하기

```
SELECT DISTINCT InvoiceNo
FROM `coral-bucksaw-482803-p2.modulabs_project.data` 
LIMIT 100
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보	결과	시각화
행	InvoiceNo	
1	541431	
2	C541433	
3	537626	
4	542237	
5	549222	
6	556201	
7	562032	
8	573511	
9	581180	
10	539318	
11	541998	

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

```
SELECT *
FROM `coral-bucksaw-482803-p2.modulabs_project.data` 
WHERE InvoiceNo LIKE 'C%' 
LIMIT 100;
```

[결과 이미지를 넣어주세요]

쿼리 결과												
작업 정보	결과	시각화	JSON	설명 세부정보	설명 그레프	StockCode	Description	Quantity	InvoiceDate	UnitPrice	CustomerID	Country
1	0541433			23166	MEDIUM CERAMIC TOP STORA...	-74215	2011-01-18 10:17:00 UTC	1	2011-03-01 15:47:00 UTC	183.75	1234	United Kingdom
2	C545329			M	Manual	-1	2011-03-01 15:47:00 UTC	1	2011-03-01 15:47:00 UTC	280.05	12352	Norway
3	C545329			M	Manual	-1	2011-03-01 15:49:00 UTC	1	2011-03-01 15:49:00 UTC	376.5	12352	Norway
4	C545330			M	Manual	-1	2011-03-01 15:49:00 UTC	1	2011-03-01 15:49:00 UTC	4.95	12352	Norway
5	C547388			22784	LANTERN CREAM GAZEBO	-3	2011-03-22 16:07:00 UTC	12	2011-03-22 16:07:00 UTC	1.25	12352	Norway
6	C547388			21914	BLUE HARMONICA IN BOX	-12	2011-03-22 16:07:00 UTC	12	2011-03-22 16:07:00 UTC	1.45	12352	Norway
7	C547388			22645	CERAMIC HEART FAIRY CAKE ...	-12	2011-03-22 16:07:00 UTC	12	2011-03-22 16:07:00 UTC	2.95	12352	Norway
8	C547388			22413	METAL SIGN TAKE IT OR LEAVE...	-6	2011-03-22 16:07:00 UTC	12	2011-03-22 16:07:00 UTC	1.49	12352	Norway
9	C547388			37448	CERAMIC CAKE DESIGN SPOTT...	-12	2011-03-22 16:07:00 UTC	12	2011-03-22 16:07:00 UTC	1.65	12352	Norway
10	C547388			84050	PINK HEART SHAPE EGG FRYIN...	-12	2011-03-22 16:07:00 UTC	12	2011-03-22 16:07:00 UTC	1.65	12352	Norway

- 구매 건 상태가 Canceled인 데이터의 비율(%) - 소수점 첫번째 자리까지

```
SELECT ROUND(SUM(CASE WHEN Quantity < 0 THEN 1 ELSE 0 END) / COUNT(*) * 100, 1)
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보 결과 시각화		
행	f0_	
1		2.2

StockCode 살펴보기

- 고유한 StockCode 의 개수를 출력하기

```
SELECT COUNT(DISTINCT StockCode)
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보 결과 시각화		
행	f0_	
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 `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY 1
ORDER BY sell_cnt DESC
LIMIT 10
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보 결과 시각화 JSON 실행 세부정보		
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 수인지 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 `coral-bucksaw-482803-p2.modulabs_project.data`
)
WHERE number_count IN(0,1);
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보	결과	시각화
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

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

```
SELECT ROUND(SUM(CASE WHEN number_count IN (0,1) THEN 1 ELSE 0 END) / COUNT(*) * 100 ,2)
FROM (
  SELECT StockCode,
    LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS number_count
  FROM `coral-bucksaw-482803-p2.modulabs_project.data`
);
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보	결과	시각화
행	f0_	
1	0.48	

- 제품과 관련되지 않은 거래 기록을 제거하기

```
DELETE FROM `coral-bucksaw-482803-p2.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 `coral-bucksaw-482803-p2.modulabs_project.data`
  )
  WHERE number_count IN(0,1)
);
```

[결과 이미지를 넣어주세요]

쿼리 결과

작업 정보 결과 실행 세부정보 실행 그래프

이 문으로 data의 행 1,915개가 삭제되었습니다.

Description 살펴보기

- 고유한 Description 별 출현 빈도를 계산하고 상위 30개를 출력하기

```
SELECT Description, COUNT(*) AS description_cnt
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY Description
ORDER BY description_cnt DESC
LIMIT 30
```

[결과 이미지를 넣어주세요]

행	Description	description_cnt
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
9	PACK OF 72 RETROSPOT CAKE ...	1062
10	SPOTTY BUNTING	1026
11	PAPER CHAIN KIT 50'S CHRIST	1013

- 서비스 관련 정보를 포함하는 행들을 제거하기

```
DELETE
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
WHERE REGEXP_CONTAINS(
    UPPER>Description),
    r'POSTAGE|ADJUSTMENT|CARRIAGE|DISCOUNT|BANK|MANUAL|WRONGLY|\?'
);
```

[결과 이미지를 넣어주세요]

```

229 DELETE FROM `coral-bucksaw-482803-p2.modulabs_project.data`
230 WHERE REGEXP_CONTAINS(
231   UPPER(`Description`),
232   r'POSTAGE|ADJUSTMENT|CARRIAGE|DISCOUNT|BANK|MANUAL|WRONGLY|\\?'
233 );
234

```

쿼리 완료됨
주문형 처리 할당량 사용 중

쿼리 결과

작업 정보 결과 실행 세부정보 실행 그래프

이 문으로 data의 행 1,314개가 삭제되었습니다.

- 대소문자를 혼합하고 있는 데이터를 대문자로 표준화 하기

```

CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.data` AS
SELECT
  * EXCEPT (`Description`),
  UPPER(`Description`) AS `Description`
FROM `coral-bucksaw-482803-p2.modulabs_project.data`;

```

[결과 이미지를 넣어주세요]

```

241 --대소문자를 혼합하고 있는 데이터를 대문자로 표준화 하기
242
243
244 CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.data` AS
245 SELECT
246   * EXCEPT (`Description`),
247   UPPER(`Description`) AS `Description`
248 FROM `coral-bucksaw-482803-p2.modulabs_project.data`;
249
250
251

```

쿼리 완료됨
주문형 처리 할당량 사용 중

쿼리 결과

작업 정보 결과 실행 세부정보 실행 그래프

이 문으로 이름이 data인 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 `coral-bucksaw-482803-p2.modulabs_project.data`;

```

[결과 이미지를 넣어주세요]

쿼리 결과

작업 정보 결과 시각화 JSON 실행 세부정보

행	min_price	max_price	avg_price
1	0.0	649.5	2.908060495763...

- 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_quantity
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
WHERE UnitPrice = 0;
```

[결과 이미지를 넣어주세요]

쿼리 결과						결과
작업 정보		결과	시각화	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

```
CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.data` AS
SELECT *
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
WHERE UnitPrice != 0;
```

[결과 이미지를 넣어주세요]

```
-- 270
271 CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.data` AS
272 SELECT *
273 FROM `coral-bucksaw-482803-p2.modulabs_project.data`
274 WHERE UnitPrice != 0;
275
```

쿼리 완료됨
주문형 처리 할당량 사용 중

쿼리 결과						결과 저장
작업 정보		결과	실행 세부정보	실행 그래프		
1	이 문으로 이름이 data인 테이블이 교체되었습니다.					

11-7. RFM 스코어

Recency

- `InvoiceDate` 컬럼을 연월일 자료형으로 변경하기

```
SELECT DATE(InvoiceDate) AS InvoiceDay,
FROM `coral-bucksaw-482803-p2.modulabs_project.data`;
```

[결과 이미지를 넣어주세요]

쿼리 결과					
작업 정보	결과	시각화	JSON	실행 세부정보	실행 그래프
행	InvoiceDay	InvoiceNo	StockCode		
1	2011-01-18	541431	23166		
2	2011-01-18	C541433	23166		
3	2010-12-07	537626	84997D		
4	2010-12-07	537626	22725		
5	2010-12-07	537626	85116		
6	2010-12-07	537626	22805		
7	2010-12-07	537626	22492		
8	2010-12-07	537626	22775		

- 가장 최근 구매 일자를 MAX() 함수로 찾아보기

```
SELECT
    DATE(MAX(InvoiceDate) OVER()) AS most_recent_date,
    DATE(InvoiceDate) AS InvoiceDay, *
FROM `coral-bucksaw-482803-p2.modulabs_project.data`;
```

[결과 이미지를 넣어주세요]

쿼리 결과											
작업 정보	결과	시각화	JSON	실행 세부정보	실행 그래프	모든 행	결과 저장	다음에서			
행	most_recent_date	InvoiceDay	InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	Description	
1	2011-12-09	2011-01-18	541431	23166	74015	2011-01-18 09:01:00 UTC	1.04	12346	United Kingdom	MEDIUM CERAMIC TOP STO...	
2	2011-12-09	2011-01-18	C541433	23166	74015	2011-01-18 10:17:00 UTC	1.04	12346	United Kingdom	MEDIUM CERAMIC TOP STO...	
3	2011-12-09	2010-12-07	537626	84997D	6	2010-12-07 14:57:00 UTC	3.75	12347	Iceland	PINK 3 PIECE POLAKOAT CUTL...	
4	2011-12-09	2010-12-07	537626	22725	4	2010-12-07 14:57:00 UTC	3.75	12347	Iceland	ALARM CLOCK BAKELIKE CH...	
5	2011-12-09	2010-12-07	537626	85116	12	2010-12-07 14:57:00 UTC	2.1	12347	Iceland	BLACK CANDLABRA T LIGHT ...	
6	2011-12-09	2010-12-07	537626	22805	12	2010-12-07 14:57:00 UTC	1.25	12347	Iceland	BLUE DRAWERS KNOB ACRYL...	
7	2011-12-09	2010-12-07	537626	22492	36	2010-12-07 14:57:00 UTC	0.65	12347	Iceland	MIN PAINT SET VINTAGE	
8	2011-12-09	2010-12-07	537626	22775	12	2010-12-07 14:57:00 UTC	1.25	12347	Iceland	PURPLE DRAWERS KNOB ACRYL...	

- 유저 별로 가장 큰 InvoiceDay를 찾아서 가장 최근 구매일로 저장하기

```
SELECT
    CustomerID,
    MAX(DATE(InvoiceDate)) AS InvoiceDay
FROM `coral-bucksaw-482803-p2.modulabs_project.data`'
GROUP BY CustomerID;
```

[결과 이미지를 넣어주세요]

작업 정보				결과	시각화	JSON
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 계산하기

```
SELECT
    CustomerID,
    EXTRACT(DAY FROM MAX(InvoiceDay) OVER () - InvoiceDay) AS recency
FROM (
    SELECT
```

```

CustomerID,
MAX(DATE(InvoiceDate)) AS InvoiceDay
FROM project_name.modulabs_project.data
GROUP BY CustomerID
);

```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보	결과	시각화
행	CustomerID	recency
1	12367	4
2	12562	8
3	12733	234
4	12824	59
5	13008	323
6	13080	179
7	13155	33
8	13391	203

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

```

CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.user_r` AS
SELECT
CustomerID,
EXTRACT(DAY FROM MAX(InvoiceDay) OVER () - InvoiceDay) AS recency
FROM (
SELECT
CustomerID,
MAX(DATE(InvoiceDate)) AS InvoiceDay
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY CustomerID
);

```

[결과 이미지를 넣어주세요]

```

317 CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.user_r` AS
318 SELECT
319   CustomerID,
320   EXTRACT(DAY FROM MAX(InvoiceDay) OVER () - InvoiceDay) AS recency
321 FROM (
322   SELECT
323     CustomerID,
324     MAX(DATE(InvoiceDate)) AS InvoiceDay
325   FROM `coral-bucksaw-482803-p2.modulabs_project.data`
326   GROUP BY CustomerID
327 );

```

✅ 쿼리 완료됨

주문형 처리 할당량 사용 중

쿼리 결과

결과 저장 ▾

작업 정보 결과 실행 세부정보 실행 그래프

ⓘ 이 문으로 이름이 user_r인 새 테이블이 생성되었습니다.

행	CustomerID	recency
1	16705	0
2	17490	0
3	18102	0
4	16626	0
5	16558	0
6	17428	0
7	13069	0
8	12985	0
9	17315	0
10	14051	0
11	12518	0
12	17364	0
13	14422	0
14	12662	0
15	13777	0

Frequency

- 고객마다 고유한 InvoiceNo의 수를 세어보기

```
SELECT
    CustomerID,
    COUNT(InvoiceNo) AS purchase_cnt
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY CustomerID ;
```

[결과 이미지를 넣어주세요]

행	CustomerID	purchase_cnt
1	12346	2
2	12347	182
3	12348	27
4	12349	72
5	12350	16
6	12352	80
7	12353	4
8	12354	58

- 각 고객 별로 구매한 아이템의 총 수량 더하기

```
SELECT
    CustomerID,
    COUNT(Quantity) AS item_cnt
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY CustomerID;
```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보	결과	시각화
행	CustomerID	item_cnt
1	12346	2
2	12347	182
3	12348	27
4	12349	72
5	12350	16
6	12352	80
7	12353	4
8	12354	58

- 전체 거래 건수 계산과 구매한 아이템의 총 수량 계산의 결과를 합쳐서 `user_rf`라는 이름의 테이블에 저장하기

```

CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.user_rf` AS

-- (1) 전체 거래 건수 계산
WITH purchase_cnt AS (
SELECT
  CustomerID,
  COUNT(InvoiceNo) AS purchase_cnt
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY CustomerID
),

-- (2) 구매한 아이템 총 수량 계산
item_cnt AS (
SELECT
  CustomerID,
  COUNT(Quantity) AS item_cnt
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
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 `coral-bucksaw-482803-p2.modulabs_project.user_r` AS ur
  ON pc.CustomerID = ur.CustomerID;

```

[결과 이미지를 넣어주세요]

```

371 -- 기존의 user_rf에 (1)과 (2)를 통합
372 SELECT
373   pc.CustomerID,
374   pc.purchase_cnt,
375   ic.item_cnt,
376   ur.reency
377 FROM purchase_cnt AS pc
378 JOIN item_cnt AS ic
379 ON pc.CustomerID = ic.CustomerID
380 JOIN `coral-bucksaw-482803-p2.modulabs_project.user_rf` AS ur
381 ON pc.CustomerID = ur.CustomerID;
382

```

쿼리 완료됨
주문형 처리 할당량 사용 중

쿼리 결과

작업 정보 결과 실행 세부정보 실행 그래프

이 문으로 이름이 user_rf인 새 테이블이 생성되었습니다.

coral-bucksaw-482803-p2 / Datasets / modulabs_project / Tables / user_rf

별명	CustomerID	purchase_cnt	item_cnt	reency
1	12791	1	1	373
2	13391	1	1	203
3	13829	1	1	359
4	15488	1	1	92
5	15668	1	1	217
6	18184	1	1	15
7	16138	1	1	368
8	13017	1	1	7
9	13307	1	1	120
10	14424	1	1	17
11	17948	1	1	147
12	12943	1	1	301
13	16953	1	1	30
14	13703	1	1	318
15	14576	1	1	372
16	16765	1	1	294
17	16078	1	1	283
18	17291	1	1	308
19	13366	1	1	50
20	13747	1	1	373
21	18141	1	1	360

Monetary

- 고객별 총 지출액 계산 (소수점 첫째 자리에서 반올림)

```

SELECT
CustomerID,
ROUND(SUM(UnitPrice), 1) AS user_total
FROM `coral-bucksaw-482803-p2.modulabs_project.data`
GROUP BY CustomerID;

```

[결과 이미지를 넣어주세요]

쿼리 결과		
작업 정보	결과	시각화
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	324.7
7	12353	24.3
8	12354	261.2

• 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 `coral-bucksaw-482803-p2.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, 1) AS user_average
FROM `coral-bucksaw-482803-p2.modulabs_project.user_rf` AS rf
LEFT JOIN (
    SELECT CustomerID, ROUND(SUM(UnitPrice), 1) AS user_total
    FROM `coral-bucksaw-482803-p2.modulabs_project.data`
    GROUP BY CustomerID
) AS ut
ON rf.CustomerID = ut.CustomerID;
```

[결과 이미지를 넣어주세요]

```
398 CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.user_rfm` AS
399 SELECT
400     rf.CustomerID AS CustomerID,
401     rf.purchase_cnt,
402     rf.item_cnt,
403     rf.recency,
404     ut.user_total,
405     ROUND(ut.user_total / rf.purchase_cnt, 1 ) AS user_average
406 FROM `coral-bucksaw-482803-p2.modulabs_project.user_rf` AS rf
407 LEFT JOIN (
408     SELECT CustomerID, ROUND(SUM(UnitPrice), 1) AS user_total
409     FROM `coral-bucksaw-482803-p2.modulabs_project.data`
410     GROUP BY CustomerID
411 ) AS ut
412 ON rf.CustomerID = ut.CustomerID;
413
414
```

● 실행 시 이 쿼리가 6.21MB를 처리합니다.

주문형 처리 할당량 사용 중

결과 저장 ·

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

행	CustomerID	purchase_cnt	item_cnt	recency	user_total	user_average
1	18141	1	1	360	3.0	3.0
2	13185	1	1	267	6.0	6.0
3	14090	1	1	324	1.1	1.1
4	17948	1	1	147	2.5	2.5
5	18133	1	1	212	0.7	0.7
6	14679	1	1	371	2.5	2.5
7	13017	1	1	7	4.3	4.3
8	13120	1	1	238	2.5	2.5
9	16148	1	1	296	1.1	1.1
10	18113	1	1	368	1.1	1.1
11	15195	1	1	2	2.8	2.8
12	16738	1	1	297	1.3	1.3
13	17763	1	1	263	1.3	1.3
14	16454	1	1	64	3.0	3.0
15	17715	1	1	200	0.8	0.8
16	16579	1	1	365	2.5	2.5
17	17443	1	1	219	1.1	1.1
18	15488	1	1	92	1.1	1.1
19	18068	1	1	289	16.9	16.9
20	18233	1	1	325	110.0	110.0
21	17331	1	1	123	10.0	10.0

RFM 통합 테이블 출력하기

- 최종 `user_rfm` 테이블을 출력하기

```
SELECT *
FROM `coral-bucksaw-482803-p2.modulabs_project.user_rfm`
```

[결과 이미지를 넣어주세요]

행	CustomerID	purchase_cnt	item_cnt	recency	user_total	user_average
1	18141	1	1	360	3.0	3.0
2	13185	1	1	267	6.0	6.0
3	14090	1	1	324	1.1	1.1
4	17948	1	1	147	2.5	2.5
5	18133	1	1	212	0.7	0.7
6	14679	1	1	371	2.5	2.5
7	13017	1	1	7	4.3	4.3
8	13120	1	1	238	2.5	2.5
9	16148	1	1	296	1.1	1.1
10	18113	1	1	368	1.1	1.1
11	15195	1	1	2	2.8	2.8
12	16738	1	1	297	1.3	1.3

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 project_name.modulabs_project.user_data AS
WITH unique_products AS (
SELECT
    CustomerID,
    COUNT(DISTINCT StockCode) AS unique_products
FROM project_name.modulabs_project.data
GROUP BY CustomerID
)
SELECT ur.*, up.* EXCEPT (CustomerID)
FROM project_name.modulabs_project.user_rfm AS ur
JOIN unique_products AS up
ON ur.CustomerID = up.CustomerID;

```

[결과 이미지를 넣어주세요]

The screenshot shows a database interface with a code editor at the top containing the SQL query for creating a table. The code includes a WITH clause for 'unique_products', a SELECT statement from 'data', and a final SELECT statement joining 'user_rfm' and 'unique_products' tables. A green checkmark icon indicates the query was successful. Below the editor is a message: '주문형 처리 할당량 사용 중' (Using ordered processing allocation). The main area shows the '쿼리 결과' (Query Results) tab selected, displaying a table with 20 rows of user data. The table has columns: 행 (Row), CustomerID, purchase_cnt, item_cnt, recency, user_total, user_average, and unique_prod... (unique_products). The data shows various customer IDs and their corresponding metrics.

행	CustomerID	purchase_cnt	item_cnt	recency	user_total	user_average	unique_prod...
1	17832	61	61	49	67.8	1.1	61
2	16024	61	61	12	70.5	1.2	60
3	12371	62	62	59	204.1	3.3	62
4	12534	62	62	130	173.2	2.8	62
5	17070	62	62	114	157.6	2.5	62
6	16644	62	62	177	95.4	1.5	62
7	13635	62	62	67	149.2	2.4	62
8	16685	62	62	61	139.8	2.3	60
9	16445	63	63	33	165.8	2.6	62
10	13221	64	64	240	326.9	5.1	60
11	14856	64	64	36	197.3	3.1	64
12	17100	65	65	18	126.3	1.9	65
13	12611	65	65	52	173.9	2.7	65
14	16795	66	66	365	145.5	2.2	64
15	15019	66	66	266	144.7	2.2	64
16	16208	66	66	44	208.5	3.2	61
17	14140	67	67	3	127.9	1.9	60
18	12762	67	67	7	132.8	2.0	61
19	16188	67	67	44	148.9	2.2	65
20	13500	67	67	23	238.3	3.6	60

The screenshot shows the BigQuery web interface with the 'user_data' table selected. The table has 20 rows of data, matching the results shown in the previous screenshot. The columns are: 행 (Row), CustomerID, purchase_cnt, item_cnt, recency, user_total, user_average, and unique_prod... (unique_products). The data is presented in a clean, tabular format with horizontal and vertical grid lines.

2. 평균 구매 주기

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

```
CREATE OR REPLACE TABLE coral-bucksaw-482803-p2.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
`coral-bucksaw-482803-p2.modulabs_project.data_backup_10min`
WHERE CustomerID IS NOT NULL
)
GROUP BY CustomerID
)

SELECT u.*, pi.* EXCEPT (CustomerID)
FROM coral-bucksaw-482803-p2.modulabs_project.user_data AS u
LEFT JOIN purchase_intervals AS pi
ON u.CustomerID = pi.CustomerID;
```

[결과 이미지를 넣어주세요]

```
477 -- (1) 구매와 구매 사이에 소요된 일수
478 SELECT
479     CustomerID,
480     DATE_DIFF(InvoiceDate, LAG(InvoiceDate) OVER (PARTITION
481 AS interval_
482     FROM
483     `coral-bucksaw-482803-p2.modulabs_project.data_backup_`  
484     WHERE CustomerID IS NOT NULL
485 )
486     GROUP BY CustomerID
487 )
488 SELECT u.*, pi.* EXCEPT (CustomerID)
489 FROM `coral-bucksaw-482803-p2.modulabs_project.user_data` AS
490 LEFT JOIN purchase_intervals AS pi
491 ON u.CustomerID = pi.CustomerID;
492
✓ 쿼리 완료됨
주문형 처리 할당량 사용 중
```

쿼리 결과

작업 정보	결과	실행 세부정보	실행 그래프

❶ 이 문으로 이름이 user_data인 테이블이 교체되었습니다.

кури 결과										결과 사용	
작성 정보		결과	시작일	JSON	실행 세부정보	실행 그레프					
행	Invoiceno	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	Description		average_interval	
1	541431	23166	74215	2010-11-18 10:01:00 UTC	1.04	12346	United Kingdom	MEDIUM CERAMIC TOP STO...		0.0	
2	2541433	23166	74215	2010-11-18 10:17:00 UTC	1.04	12346	United Kingdom	MEDIUM CERAMIC TOP STO...		0.0	
3	537626	21171	12	2010-12-07 14:57:00 UTC	1.45	12347	Iceland	BATHROOM METAL SIGN		2.0	
4	537626	848970	6	2010-12-07 14:57:00 UTC	3.75	12347	Iceland	BLUE 3 PIECE POLKA DOT CUTL...		2.0	
5	537626	22212	6	2010-12-07 14:57:00 UTC	2.1	12347	Iceland	FOUR HORN WHITE LOVEBIRD		2.0	
6	537626	851678	30	2010-12-07 14:57:00 UTC	1.25	12347	Iceland	BLACK GRAND BARBOUR PHOT...		2.0	
7	537626	22729	4	2010-12-07 14:57:00 UTC	3.75	12347	Iceland	ALARM CLOCK BAXIELKE ORA...		2.0	
8	537626	23497	4	2010-12-07 14:57:00 UTC	4.25	12347	Iceland	SET OF 2 TINS VINTAGE BATH...		2.0	
9	537626	23805	12	2010-12-07 14:57:00 UTC	1.25	12347	Iceland	BLUE DRAWER KNOK ACRYLIC...		2.0	
10	537626	22726	4	2010-12-07 14:57:00 UTC	3.75	12347	Iceland	ALARM CLOCK BAXIELKE GREEN		2.0	

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 `coral-bucksaw-482803-p2.modulabs_project.user_data` AS

WITH TransactionInfo AS (
    SELECT
        CustomerID,
        COUNT(8) AS total_transactions,
        SUM(CASE WHEN Quantity < 0 THEN 1 ELSE 0 END ) AS cancel_frequency
    FROM `coral-bucksaw-482803-p2.modulabs_project.data_backup_10min`
    WHERE CustomerID IS NOT NULL
    GROUP BY CustomerID
)

SELECT u.* , t.* EXCEPT(CustomerID), SAFE_DIVIDE(t.cancel_frequency, t.total_transactions ) AS cancel_rate
FROM `coral-bucksaw-482803-p2.modulabs_project.user_data` AS u
LEFT JOIN TransactionInfo AS t
ON u.CustomerID = t.CustomerID;
```

[결과 이미지를 넣어주세요]

```
503 -- 퀴리 쿼리, 바로 실행
504
505 CREATE OR REPLACE TABLE `coral-bucksaw-482803-p2.modulabs_project.user_data` AS
506
507 WITH TransactionInfo AS (
508     SELECT
509         CustomerID,
510         COUNT(8) AS total_transactions,
511         SUM(CASE WHEN Quantity < 0 THEN 1 ELSE 0 END ) AS cancel_frequency
512     FROM `coral-bucksaw-482803-p2.modulabs_project.data_backup_10min`
513     WHERE CustomerID IS NOT NULL
514     GROUP BY CustomerID
515 )
516
517
518     SELECT u.* , t.* EXCEPT(CustomerID), SAFE_DIVIDE(t.cancel_frequency, t.total_transactions)
519     FROM `coral-bucksaw-482803-p2.modulabs_project.user_data` AS u
520     LEFT JOIN TransactionInfo AS t
521     ON u.CustomerID = t.CustomerID;
522
523
524 ! CREATE TABLE has columns with duplicate name total_transactions at [509:1]
525
526 주문형 처리 할당량 사용 중
```

쿼리 결과

작업 정보 결과 실행 세부정보 실행 그래프

이 문으로 이름이 user_data인 테이블이 교체되었습니다.

쿼리 결과										DB 결과 차량			
작업 정보		결과		시작일		JSON		실행 세부정보		실행 그레프			
#	Invoicedate	StockCode	Quantity	Invoicedate	UnitPrice	CustomerID	Country	Description	average_interval	total_transactions	cancel_frequency	cancel_rate	
1	541431	23166	74215	2011-01-18 01:00:00	1.04	12346	United Kingdom	MEDIUM CERAMIC TOP STO...	0.0	2	1	0.0	
2	C514345	23166	74215	2011-01-18 17:00:00	1.04	12346	United Kingdom	MEDIUM CERAMIC TOP STO...	0.0	2	1	0.0	
3	S37626	84997C	6	2010-12-17 04:57:00	3.75	12347	Iceland	BLUE 5 PIECE POLAKOUD CUTL...	2.0	182	0	0.0	
4	S37626	22728	6	2010-12-17 04:57:00	3.75	12347	Iceland	ALARM CLOCK BAKELIKE PINK	2.0	182	0	0.0	
5	S37626	22212	6	2010-12-17 04:57:00	2.1	12347	Iceland	FOUR HOOF WHITE LOVEBIRDS	2.0	182	0	0.0	
6	S37626	21791	12	2010-12-17 04:57:00	1.65	12347	Iceland	RED TOADSTOOL, LED NIGHT L...	2.0	182	0	0.0	

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

```
SELECT *  
FROM `coral-bucksaw-482803-p2.modulabs_project.user_data`
```

[결과 이미지를 넣어주세요]

쿼리 결과									
작업 정보		JSON		설정 세부정보		설명 그레프			
번호	InvoiceNo	StockCode	Quantity	InvoiceDate	UnitPrice	CustomerID	Country	Description	average_interval
1	541431	23166	74215	2011-01-18 10:01:00 UTC	1.04	1234	United Kingdom	MEDIUM CERAMIC TOP STO...	0.0
2	C541432	23166	-74215	2011-01-18 10:17:00 UTC	1.04	1234	United Kingdom	MEDIUM CERAMIC TOP STO...	0.0
3	537626	849970	6	2010-12-07 14:57:00 UTC	3.75	12347	Iceland	BLUE 3 PIECE POLKA DOT OUTL...	2.0
4	537626	22728	4	2010-12-07 14:57:00 UTC	3.75	12347	Iceland	ALARM CLOCK BAKELIKE PINK	2.0
5	537626	22212	6	2010-12-07 14:57:00 UTC	2.1	12347	Iceland	FOUR HOOK WHITE LOVEBIRDS	2.0
6	537626	21791	12	2010-12-07 14:57:00 UTC	1.65	12347	Iceland	RED TOASTSTOOL LED NIGHT LI...	2.0

중간에 데이터 한번 날려 먹어서 진행이 안되는 것일 수도 있어요... 마지막에 안맞네요ㅠ

회고

[회고 내용을 작성해주세요]

Keep :

Problem : 백업테이블 없이 CREATE OR REPLACE TABLE 하다가 테이블 이름 잘 못 적어서 기본 데이터인 data 테이블을 날려먹음

Try : 10분전 데이터 새 테이블로 복사해서 시도