## 컬럼 별 누락된 값의 비율 계산

그럼, 각 컬럼 별 누락된 값의 비율을 계산해보겠습니다.

각 컬럼에 대해서 누락 값을 계산한 후, 계산된 누락 값을 UNION ALL을 통해 합쳐주세요.

```
16 SELECT
      'InvoiceNo' AS column_name,
17
18
     ROUND (
19
       SUM(CASE WHEN InvoiceNo IS NULL THEN 1 ELSE @ END)
20
       / COUNT(*) * 100
21
     ) AS missing_percentage
22
23 FROM praxis-beacon-464902-t2.modulabs_project.data
24
25 UNION ALL
26
27 SELECT
28 'StockCode' AS column_name,
29
     ROUND (
       SUM(CASE WHEN StockCode IS NULL THEN 1 ELSE @ END)
30
        / COUNT(*) * 100
31
   , 2
) AS missing_percentage
32
33
34 FROM praxis-beacon-464902-t2.modulabs_project.data
35
36 UNION ALL
37
38 SELECT
39
      'Description' AS column_name,
40
      ROUND (
       SUM(CASE WHEN Description IS NULL THEN 1 ELSE @ END)
41
        / COUNT(*) * 100
42
43
     ) AS missing_percentage
44
45 FROM praxis-beacon-464902-t2.modulabs_project.data
46
47 UNION ALL
48
49 SELECT
      'Quantity'
50
                    AS column_name,
      ROUND (
51
52
       SUM(CASE WHEN Quantity IS NULL THEN 1 ELSE @ END)
53
       / COUNT(*) * 100
```

```
54 , 2
55 ) AS missing_percentage
56 FROM praxis-beacon-464902-t2.modulabs_project.data
57
58 UNION ALL
59
60 SELECT
61
     'InvoiceDate' AS column_name,
62
     ROUND(
63
       SUM(CASE WHEN InvoiceDate IS NULL THEN 1 ELSE 0 END)
       / COUNT(*) * 100
64
    , 2
) AS missing_percentage
65
66
67 FROM praxis-beacon-464902-t2.modulabs_project.data
68
69 UNION ALL
70
71 SELECT
72
     'UnitPrice' AS column_name,
73
      ROUND (
74
       SUM(CASE WHEN UnitPrice IS NULL THEN 1 ELSE 0 END)
75
       / COUNT(*) * 100
76
77
    ) AS missing_percentage
78 FROM praxis-beacon-464902-t2.modulabs_project.data
79
80 UNION ALL
81
82 SELECT
83
     'CustomerID'
                    AS column_name,
84
     ROUND (
85
       SUM(CASE WHEN CustomerID IS NULL THEN 1 ELSE 0 END)
86
       / COUNT(*) * 100
87
    ) AS missing_percentage
89 FROM praxis-beacon-464902-t2.modulabs_project.data
90
91 UNION ALL
92
93 SELECT
                     AS column_name,
 94
      'Country'
 95
     ROUND(
 96
       SUM(CASE WHEN Country IS NULL THEN 1 ELSE 0 END)
 97
        / COUNT(*) * 100
98
       , 2
99
    ) AS missing_percentage
100 FROM praxis-beacon-464902-t2.modulabs_project.data;
191
```

# ← 쿼리 결과

작업 정	보	결과	차트	JSON	실행 세부정보	1
행 <i>//</i>	column_	name 🕶		// missir	ng_percentage • //	
1	UnitPrice	!			0.0	
2	Country				0.0	
3	Descripti	on			0.27	
4	Custome	rID			24.93	
5	InvoiceN	0			0.0	
6	Quantity				0.0	
7	InvoiceD	ate			0.0	
8	StockCo	de			0.0	

# 결측치 처리 전략



# 결측치 처리

```
DELETE FROM praxis-beacon-464902-t2.modulabs_project.data2
109
      WHERE
110
111
        Description IS NULL OR
112
113
        CustomerID IS NULL;
114
 쿼리 결과
작업 정보
                                    실행 그래프
            결과
                     실행 세부정보
      이 문으로 data2의 행 135,080개가 삭제되었습니다.
 0
```

# 11-5. 데이터 전처리(2): 중복값 처리

# 중복값 확인

```
116
     SELECT
       COUNT(*) AS duplicate_group_count
117
118 FROM (
119
       SELECT
         COUNT(*) AS cnt
120
121
       FROM
         praxis-beacon-464902-t2.modulabs_project.data2
122
       GROUP BY
123
         InvoiceNo,
124
125
         StockCode,
126
         Description,
127
         Quantity,
128
         InvoiceDate,
129
         UnitPrice,
130
         CustomerID,
131
         Country
132
       HAVING
133
         cnt > 1
134 ) t;
        쿼리 결과
  작업 정보
                결과
                         차트
         duplicate_group_... //
     1
                    4837
```

## 중복값 처리

```
135
136
     CREATE OR REPLACE TABLE praxis-beacon-464902-t2.modulabs_project.data2 AS
137
     DISTINCT *
138
     FROM
139
140
       praxis-beacon-464902-t2.modulabs_project.data2;
      쿼리 결과
 -
                     실행 세부정보
 작업 정보
             결과
                                    실행 그래프
       이 문으로 이름이 data2인 테이블이 교체되었습니다.
  0
142
      SELECT
143
      COUNT(*) AS remaining_rows
144
   FROM
145 praxis-beacon-464902-t2.modulabs_project.data2;
 ← 쿼리결과
 작업 정보
             결과
                     大
       remaining_rows -
    1
               401604
```

# 11-6. 데이터 전처리(3): 오류값 처리

# InvoiceNo 살펴보기

고유(unique)한 InvoiceNo의 개수를 출력해 보세요.

```
147 SELECT
148 COUNT(DISTINCT InvoiceNo) AS unique_invoiceno_count
149 FROM
150 praxis-beacon-464902-t2.modulabs_project.data2;
```

# Q. 유니크한 InvoiceNo는 몇 개인가요?



이번에는 고유한 InvoiceNo를 100개를 출력해 보세요. 아래와 같은 결과가 나와야 합니다.

```
152 WITH unique_invoices AS (
153
      SELECT
        DISTINCT InvoiceNo
154
155
       FROM
        praxis-beacon-464902-t2.modulabs_project.data2
156
157
158 SELECT
159 InvoiceNo
160 FROM
161
     unique_invoices
162 LIMIT
163 100;
```



출력된 결과를 봤을 때, 특이한 점이 있는 값들이 있나요?

데이터를 자세히 살펴보면 몇 개의 InvoiceNo는 'C'로 시작한다는 것을 알 수 있습니다. 데이터 설명에도 나와 있던 것처럼, InvoiceNo가 'C'로 시작한다면 취소한 거래입니다.

고객 행동과 제품 선호도에 대한 이해를 높이기 위해서는 취소된 거래들도 고려해야 합니다.

먼저, InvoiceNo가 'C'로 시작하는 행을 필터링하여 취소된 거래들만 살펴봅시다. 이후에는 이 행들을 분석하여 공통적인 특성이나 패턴이 있는지 파악할 것입니다.

InvoiceNo가 'C'로 시작하는 행을 필터링할 수 있는 쿼리문을 작성해 주세요.

100행까지만 출력해 주세요. 아래와 같이 결과가 나와야 합니다.



데이터에 특이한 경향성이 보이나요?

- Q. 취소된 거래 건들은 어떠한 특징이나 경향성이 있나요?
  - Quantity가 음수입니다.

구매 건 상태가 Canceled 인 데이터의 비율(%)은 어떻게 되나요? 이를 계산할 수 있는 쿼리 문을 작성하고, 취소 비율을 소수점 첫번째 자리까지 구해 주세요

```
173 SELECT
174
       ROUND (
175
         SUM (CASE
               WHEN STARTS_WITH(InvoiceNo, 'C') THEN 1
176
177
               ELSE 0
178
             END)
        / COUNT(*) * 100
179
180
       , 1) AS canceled_percentage
181
182 praxis-beacon-464902-t2.modulabs_project.data2;
       쿼리 결과
 \leftarrow
 작업 정보
               결과
        canceled_percent...
    1
                    2.2
```

데이터를 자세히 보면 InvoiceNo 컬럼에서 'C'가 붙은 거래 취소 건들은 Quantity가 음수인 것을 알 수 있습니다.

취소를 많이 한 제품들의 가격대가 높았는지, 또는 거래 지역이 특정 지역에 몰려 있는지도 살펴보았지만, 그런 경향성은 크게 보이지 않습니다. (Country는 취소 여부와 상관없이 United Kingdom에 밀집되어 있습니다.)

그럼 InvoiceNo에 따른 추가 처리를 해야 할까요?

실제 분석 과정에서는 분석가의 주관적 판단에 의해 전처리 여부와 그 방식을 결정 짓습니다. 위의 경우 프로젝트의 초기 목표가 고객들의 구매 최신성, 구매 빈도, 구매 금액에 따라 세그멘테이션하는 것이었기 때문에, 고객의 취소 패턴을 이해하는 것도 중요할 것입니다. 가령 취소된 거래에 공통점이 있는지 살펴볼 수도 있을 것이고, 추천될 가능성이 높은 제품을 찾아서 문제를 진단할 수도 있습니다.

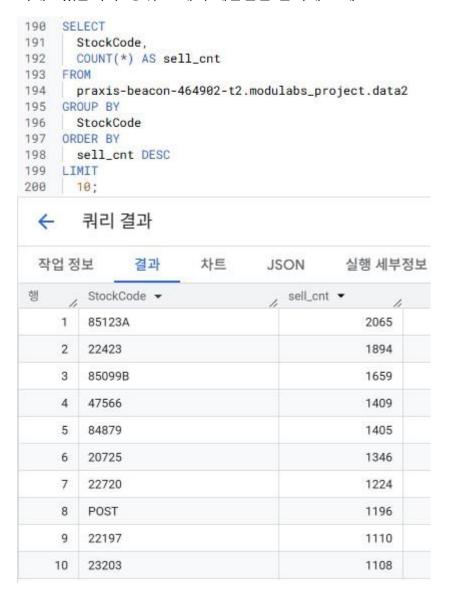
따라서 취소된 거래 데이터는 유지하되, 명확하게 표시하여 추가 분석을 용이하게 만들어주는 것도 전략 중 하나입니다.

## StockCode 살펴보기

이번엔 StockCode를 더 깊이 있게 살펴볼 차례입니다. 우선 고유한 StockCode의 개수를 출력해보겠습니다.



이번에는 어떤 제품이 가장 많이 판매되었는지 보기 위하여 StockCode 별 **등장 빈도**를 출력해보겠습니다. 상위 10개의 제품들을 출력해 보세요.



- Q. StockCode 별 등장 빈도를 출력했을 때, 일정한 경향성이 보이나요? 이상치는 있었나요?
  - 대부분의 StockCode들은 5-6 자리 숫자입니다. StockCode 중 'POST'는 이상치 같아 보입니다.

가장 판매가 많이 일어난 상위 10개의 제품 코드를 자세히 살펴보면, 고객들이 자주 구매하는 인기 제품군이나 카테고리에 대한 인사이트도 얻을 수가 있습니다.

제품 코드는 대부분 5~6자리의 숫자와 문자 조합으로 구성되어 있는 반면, 'POST'와 같은 몇 가지 이상한 코드도 있습니다. 이러한 이상 현상은 실제 제품보다는 서비스나 배송비 같은 형태를 코드로 남긴 것일 수도 있습니다. 현재 진행하고 있는 프로젝트는 고객들의 제품 구매에 초점이 맞춰져 있기 때문에 이런 값들은 제거하는 것이 좋을 것 같습니다.

'POST'와 같은 이상치들이 몇 개나 있는지 확인하기 위하여 StockCode의 문자열 내 숫자의 길이를 살펴보겠습니다. 예를 들어 '22423'와 '85123A'는 모두 다 숫자가 5개씩 포함되어 있는 문자열이지만 'POST'와 같은 비정상적인 항목들은 숫자가 0개 포함되어 있습니다.

이처럼 각 StockCode 내 숫자의 개수를 살펴보면서 값의 특성에 대한 통찰력을 얻을 수 있습니다.

## StockCode의 **문자열 내 숫자의 길이**를 구해 봅시다.

```
WITH UniqueStockCodes AS (

SELECT DISTINCT StockCode
FROM praxis-beacon-464902-t2.modulabs_project.data2

05 )

SELECT
LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS number_count,
COUNT(*) AS stock_cnt
FROM UniqueStockCodes
GROUP BY number_count
ORDER BY stock_cnt DESC;
```

위의 코드를 더 자세히 살펴보도록 하겠습니다.

우선 REGEXP\_REPLACE 라는 함수는 텍스트를 처리하는 정규 표현식(Regular Expression) 중 하나입니다. 'REGEXP'는 정규 표현식를 의미하며, 'REPLACE'는 텍스트를 대체한다는 의미로, REGEXP\_REPLACE 함수는 특정 조건에 부합한 텍스트를 다른 텍스트로 대체하는 함수입니다.

LENGTH 함수 내부에 있는 코드는 REGEXP\_REPLACE(StockCode, r'[0-9]', ")라고 작성되어 있습니다. 이 코드는 StockCode 컬럼에 있던 값 중에서 0부터 9 사이의 숫자([0-9])를 비어 있는 값(")으로 대체하는 코드입니다. 이 코드를 통해 숫자를 제외한 문자만 남게 됩니다. 이후에 LENGTH 함수로 감싸주어서, 각 StockCode에 문자가 몇 자리 수인지를 세어주었습니다.

최종적으로 LENGTH(StockCode) - LENGTH(REGEXP\_REPLACE(StockCode, r'[0-9]', '') 연산을 통하여 StockCode 안에 있는 숫자의 수를 세어준 후, 이를 number\_count라는 이름의 컬럼으로 저장해 주었습니다.

쿼리문을 실행시키면 아래와 같은 결과가 나옵니다.

<b>←</b>	쿼리 결과			
작업 정.	보 결과	t	上	JSON
행 //	number_count •	,	stock_	ent ▼
1		5		3676
2		0		7
3		1		1

출력 결과를 보면, 8개를 제외하곤 StockCode에 **5개의 숫자들이 포함**되어 있는 것을 알 수 있습니다. 숫자가 0 개인 코드는 7개, 숫자가 1개인 코드는 1개입니다.

# 숫자가 0~1개인 값들에는 어떤 코드들이 들어가 있는지를 확인해 봅시다.

```
213 SELECT DISTINCT StockCode, number_count
214 FROM (
215 | SELECT StockCode,
216 | LENGTH(StockCode) - LENGTH(REGEXP_REPLACE(StockCode, r'[0-9]', '')) AS number_count
217 | FROM praxis-beacon-464902-t2.modulabs_project.data2
218 )
219 WHERE number_count BETWEEN @ AND 1;
```

# Q. 숫자가 0~1개인 값들에는 어떤 코드들이 들어가 있나요?

- POST, D, C2, M, BANK CHARGES, PADS, DOT, CRUK

<b>←</b>	쿼리	결과			
작업 정	보	결과	차트	JSOI	N 실행 세부정
행 //	Stock	kCode ▼		/ nu	umber_count • //
1	POST	Γ			0
2	М				0
3	C2				1
4	D				0
5	BANI	K CHARGES			0
6	PADS	S			0
7	DOT				0
8	CRU	K			0

해당 코드 값들을 가지고 있는 데이터 수는 전체 데이터 수 대비 몇 퍼센트일까요? 소수점 두번째 자리까지 구해 주세요.

```
221 SELECT
       ROUND(
222
223
         COUNTIF(
224
           -- StockCode 내 숫자 개수 계산 후 0~1 사이인지 판별
225
           LENGTH(StockCode)

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

226
227
           BETWEEN 0 AND 1
228
229
        / COUNT(*) * 100
       , 2) AS percent_zero_or_one_digits
230
231
     FROM
     praxis-beacon-464902-t2.modulabs_project.data2;
232
```

# Q. 해당 코드 값들을 가지고 있는 데이터 수는 전체 데이터 수 대비 몇 퍼센트인가요?

- 0.48%



분석에 따르면 전체 데이터 중 매우 작은 비율인 0.48%의 데이터가 일반적인 형식에서 벗어난 값을 가지고 있었습니다.

이 코드들은 'BANK CHARGES, POST' 등 제품과 관련되지 않은 거래 기록으로 보입니다. 진행하고 있는 프로젝트의 목표는 고객들의 '제품 구매'에 기반하여 세그먼테이션을 하는 것이므로, 이런 StockCode가 포함된 기록은 데이터셋에서 제외하도록 하겠습니다.

## 제품과 관련되지 않은 거래 기록을 제거하는 쿼리문을 작성해 주세요.

```
DELETE FROM praxis-beacon-464902-t2.modulabs_project.data2

WHERE StockCode IN (

SELECT DISTINCT StockCode

FROM (

SELECT StockCode,

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

FROM praxis-beacon-464902-t2.modulabs_project.data2)

where number_count between 0 and 1);
```

# Description 살펴보기

먼저 데이터셋에서 고유한 Description 별 출현 빈도를 계산하고 상위 30개를 출력해 주세요.

```
243 SELECT
244 Description,
245 COUNT(*) AS description_count
246 FROM
247 praxis-beacon-464902-t2.modulabs_project.data2
248 GROUP BY
249 Description
250 ORDER BY
251 description_count DESC
252 LIMIT
253 30;
```

출력 결과를 보니 제품군들이 대부분 가정용품, 주방용품, 점심 도시락, 장식품과 관련된 것들임을 알 수 있습니다.

흥미롭게도 모든 Description이 대문자로 되어 있는데, 이는 데이터베이스에 제품 설명을 입력할 때 사용되는 표준화된 형식일 수도 있습니다. 그러나 혹시라도 대문자와 소문자가 혼합된 스타일로 입력된 설명이 있는지 확인해 보는 것이 현명할 것입니다.

# 대소문자가 혼합된 Description이 있는지 확인해 봅시다.

```
255 SELECT DISTINCT Description
256 FROM praxis-beacon-464902-t2.modulabs_project.data2
257 WHERE REGEXP_CONTAINS(Description, r'[a-z]');
```

위의 코드를 더 자세히 살펴보도록 하겠습니다.

WHERE 절의 REGEXP\_CONTAINS 함수는 특정 패턴이 문자열에 포함되어 있는지 여부를 확인하는 데에 유용한함수입니다. 특정 패턴이 문자열에 포함되어 있으면 True를, 포함되어 있지 않으면 False를 반환합니다. REGEXP\_CONTAINS(Description, r'[a-z]') 코드는 Description 컬럼에 있는 문자열에서 소문자 알파벳([a-z])이 포함되어 있는지를 확인하는 조건문입니다. 만약 r'[a-z]' 대신 r'[A-Z]' 를 사용했다면 대문자 알파벳이 포함되어 있는지를 확인하는 조건문이 됩니다.

결과는 아래와 같이 나옵니다. 빅쿼리에서 실행 결과의 스크롤을 내리면 총 19개의 Description이 대소문자를 혼합하고 있다는 것을 알 수 있습니다.

"3 TRADITIONAL BISCUIT CUTTERS SET'는 모두 대문자로 작성된 거 아냐?"라고 생각할 수 있지만 TRADITIONAL'에서 'I'은 대문자 I(아이)가 아니라 소문자 I(엘)입니다.

# ← 쿼리결과

작업 정	정보 결과	차트	JSON	실행 세부정보	
행	Description •	_		description_count	7/
1	WHITE HANGII	NG HEART T-LI	GHT HOLDER	2058	
2	REGENCY CAK	ESTAND 3 TIE	1894		
3	JUMBO BAG R	ED RETROSPO	1659	(	
4	PARTY BUNTIN	IG	1409	ĺ	
5	ASSORTED CO	1405	į		
6	LUNCH BAG RI	ED RETROSPO	Т	1345	i
7	SET OF 3 CAKE	TINS PANTRY	1224		
8	LUNCH BAG B	LACK SKULL.		1099	
9	PACK OF 72 RE	TROSPOT CA	KE CASES	1062	
10	SPOTTY BUNT	ING		1026	i
11	PAPER CHAIN	KIT 50'S CHRIS	STMAS	1013	
12	LUNCH BAG SE	PACEBOY DESI	GN	1006	1
13	LUNCH BAG C	ARS BLUE		1000	)
14	HEART OF WIC	KER SMALL		990	1
15	NATURAL SLA	TE HEART CHA	ALKBOARD	989	
16	JAM MAKING	SET WITH JAR	S	966	
17	LUNCH BAG PI	NK POLKADO	Γ	961	
18	LUNCH BAG SI	JKI DESIGN		932	,

출력 결과를 보면 사이즈(cm)나 무게(g) 등의 단위를 나타내는 설명이 포함되어 있고, 'Next Day Carriage'나 'High Resolution Image'와 같은 일부 항목들처럼 실제 제품에 대한 Description이 아닌 것도 있는 것을 알 수 있습니다. 이 데이터들은 실제 제품 정보와 관련이 없어 보입니다. 아마도 다른 유형의 정보나 서비스 세부사항을 나타내는 것일수도 있습니다. 이를 처리하기 위한 몇가지 전략들을 생각해볼 수 있습니다.

- 1. 'Next Day Carriage'와 'High Resolution Image'와 같은 서비스 관련 정보를 포함하는 행들을 제거합니다.
- 2. 대소문자를 혼합해서 사용하는 경우, 대문자로 표준화하여 데이터셋 전체에서 일관성을 유지할 수 있습니다. 이는 대소문자에 의한 중복 항목의 가능성을 줄이는 데에도 도움이 될 것입니다.

이러한 전략을 선택함으로써, 데이터셋의 품질을 향상시킬 수 있고 프로젝트의 분석 단계에 더 적합한 데이터셋을 완성할 수 있습니다.

우선 **서비스 관련 정보를 포함하는 행들을 제거**하는 쿼리문을 작성해 보세요.

```
259 DELETE
260 FROM praxis-beacon-464902-t2.modulabs_project.data2
261 WHERE
262 -- Description에 서비스 관련 키워드가 포함된 행 제거
263 REGEXP_CONTAINS(
LOWER(Description),
r'service|delivery|postage|shipping|charges'
266 );
```

# 이번에는 **대소문자를 혼합하고 있는 데이터를 대문자로 표준화**하는 쿼리문을 작성해 보세요.

```
269 CREATE OR REPLACE TABLE praxis-beacon-464902-t2.modulabs_project.data2 AS
270 SELECT
271 * EXCEPT (Description),
272 UPPER(Description) AS Description
273 FROM
274 praxis-beacon-464902-t2.modulabs_project.data2;
```

SELECT 절에 EXCEPT (Description)이라는 구문은 처음 보셨을 겁니다. \* EXCEPT (Description)는 테이블의 모든 컬럼들을 가져오되, Description은 제외하라는 의미입니다. 포함하려고 하는 컬럼명을 일일이 쓰는 것보다 제거하고자 하는 몇 개 컬럼을 적을 때 자주 활용합니다.

Quantity와 InvoiceDate에 대한 데이터를 분석했을 때에는 크게 문제가 되는 이상 데이터가 없어서 넘어가겠습니다. 데이터를 분석할 때 전처리하면 좋을 것 같은 데이터가 나타난다면, 분석가/과학자의 판단에 따라 처리 방법을 선택할 수 있습니다.

## UnitPrice 살펴보기

이번에는 UnitPrice에서 이상치를 찾아봅시다. 최솟값, 최댓값, 평균 데이터를 확인해 봄으로써, 단위 가격의 요약 통계량을 살펴보겠습니다.

※ 단위 가격(UnitPrice)이란?

상품 1개당 가격을 말합니다. 즉 "단위 당 가격"을 의미합니다.

## UnitPrice의 **최솟값, 최댓값, 평균**을 구해 보세요.

```
277 SELECT
278 MIN(UnitPrice) AS min_unitprice,
279 MAX(UnitPrice) AS max_unitprice,
280 ROUND(AVG(UnitPrice), 2) AS avg_unitprice
281 FROM
282 praxis-beacon-464902-t2.modulabs_project.data2;
```



# Q. UnitPrice 데이터에 특이한 값이 있나요?

- UnitPrice가 0인 데이터가 있습니다.

단위 가격의 요약 통계량을 보면 최소 단위 가격(min\_price)이 0인 것을 알 수 있습니다. 이는 단가가 0원인 데이터가 존재한다는 것을 의미하고, 이는 이 제품이 무료 제품이거나 데이터 오류일 수도 있다는 의미입니다.

단가가 0원인 거래의 성격을 제대로 이해하기 위해서는 데이터를 더 자세히 살펴볼 필요가 있습니다. 단가가 0 원인 제품을 상세하게 분석해보면서 특정한 패턴이 있는지 살펴보겠습니다.

# 단가가 0원인 거래의 개수, 구매 수량(Quantity)의 최솟값, 최댓값, 평균을 구해 보세요.

```
284 SELECT
       COUNT(*)
285
                                           AS zero_price_count,
286
       MIN(Quantity)
                                           AS min_quantity,
287
       MAX(Quantity)
                                           AS max_quantity,
288
       ROUND(AVG(Quantity), 2)
                                           AS avg_quantity
289
       praxis-beacon-464902-t2.modulabs_project.data2
290
291
    WHERE
292 UnitPrice = 0;
```

# 쿼리 결과

작인	법 정보	<b>결과</b>	,	가트 JSON		실행 세부정보	실행 그래프
행	1,	zero_price_count	-/	min_quantity 🕶	1	max_quantity ▼ //	avg_quantity •
	1	3	33		1	12540	420.52

UnitPrice가 0인 행의 수는 33개로 비교적 적음을 알 수 있습니다. 구매 수량(Quantity)은 최소 1개부터 최대 12,540개에 이르기까지 굉장히 큰 편차를 가집니다.

데이터의 수가 적은 걸 보니 무료 제품이라기보다 데이터 오류일 가능성이 더 높을 것 같습니다. 그래서 **이 데이터** (UnitPrice = 0)를 제거하고 일관된 데이터셋을 유지하도록 하겠습니다.

```
CREATE OR REPLACE TABLE praxis-beacon-464902-t2.modulabs_project.data2 AS
SELECT *
Praxis-beacon-464902-t2.modulabs_project.data2
WHERE
UnitPrice != 0;
```

#### 11-7. RFM 스코어

이번 시간에는 본격적으로 RFM 분석을 해보도록 하겠습니다. 다시 한번 RFM 분석에 대해 배운 내용을 복습해 봅시다.

RFM분석은 구매 최신성(Recency), 구매 빈도(Frequency), 구매 가치(Monetary)에 따라 고객들을 여러 그룹으로 나누는 세그먼테이션(segmentation) 방법이라고 하였습니다.

- Recency는 고객이 마지막으로 구매한 시점을 나타냅니다. 최근에 구매한 고객들은 더 자주 구매할 가능성이 높기 때문에, 최신성 점수가 높은지를 고려합니다.
- Frequency는 특정 기간 동안 고객이 얼마나 자주 우리의 제품이나 서비스를 구매하는지를 나타냅니다. 빈번하게 구매하는 고객은 충성도가 높은 고객일 확률이 높기 때문에, 빈도 점수가 높은지를 고려합니다.
- Monetary는 고객이 지출한 총 금액을 말합니다. 많은 금액을 지불한 고객일수록 더 가치가 높은 충성 고객일 수 있습니다. 앞으로도 우리의 제품과 사이트에 많은 돈을 지불할 수 있는 고객이므로, 가치 점 수가 높은지를 함께 고려합니다.

#### Recency

Recency 단계에서는 고객이 얼마나 최근에 구매를 했는지에 중점을 둡니다. 그러므로 '마지막 구매일로부터 현재까지 경과한 일수'를 계산해야 합니다. 낮은 값일수록 고객이 최근에 구매를 했음을 의미하며, 제품이나 서비스에 더 관심을 보인다고 예측할 수 있습니다. Recency를 통해 오랜 시간 동안 구매를 하지 않았던 고객을 발견하고, 다시 제품과 서비스로 불러들이기 위한 마케팅 전략을 맞춤화해볼 수도 있습니다.

우선 InvoiceDate를 '2010-12-01 08:26:00'와 같은 'YYYY-MM-DD HH:MM:SS' 형태에서 'YYYY-MM-DD' 형태로 날짜에 해당하는 부분만 남겨놓고 싶습니다. 이를 위해 DATE 함수를 활용하여 InvoiceDate 컬럼을 **연월일 자료 형으로 변경**해 주세요.

```
301 SELECT
302 DATE(InvoiceDate) AS InvoiceDay,
303 *
304 FROM
305 praxis-beacon-464902-t2.modulabs_project.data2;
```

실제 회사에서 다루는 데이터라면 오늘 날짜를 기준으로 최종 구매일이 몇 일 지났는지 계산하겠지만, 여러분이 다루고 있는 데이터는 2010년~2011년 사이의 데이터이므로, 최종 구매일로부터 꽤 오랜 시간이 지난 데이터라는 특성이 있습니다. 그래서 이번 프로젝트에서는 모든 고객들을 통틀어 가장 최근 구매 일자를 기준으로 Recency를 구하려고 합니다.

우선 가장 최근 구매 일자를 MAX() 함수로 찾아보겠습니다.

```
307 SELECT
308
            MAX(InvoiceDate) OVER()
                                                             AS most_recent_date,
309
            DATE(InvoiceDate)
                                                              AS InvoiceDay,
310
311
        FROM
312
           praxis-beacon-464902-t2.modulabs_project.data2;
쿼리 결과
                                                                                               ③ 결과 저장 ▼
                                                                                                                 ₩ 다음에서 열기 ▼
작업 정보
             결과
                      차트
                                JSON
                                           실행 세부정보
                                                            실행 그래프
                                                 / InvoiceNo ▼
 ____ most_recent_date ▼
                                  InvoiceDay ▼
                                                                                StockCode •
                                                                                                            Quantity -
                                                                                                                             InvoiceDate ▼
  1 2011-12-09 12:50:00 UTC
                                  2011-10-13
                                                    571034
                                                                                23094
                                                                                                                             2011-10-13 12:4
  2 2011-12-09 12:50:00 UTC
                                  2011-12-02
                                                    C580165
                                                                                22826
                                                                                                                             2011-12-02 11:2
      2011-12-09 12:50:00 UTC
                                  2011-08-19
  4 2011-12-09 12:50:00 UTC
                                  2011-06-17
                                                    557147
                                                                                47518F
                                                                                                                             2011-06-17 10:5
                                                                                                                        24
      2011-12-09 12:50:00 UTC
                                  2011-03-03
                                                    545475
                                                                                21915
                                                                                                                             2011-03-03 10:5
  6 2011-12-09 12:50:00 UTC
                                  2011-08-18
                                                    563614
                                                                                23203
                                                                                                                        100
                                                                                                                             2011-08-18 08:5
      2011-12-09 12:50:00 UTC
                                  2011-08-18
                                                    563614
                                                                                23343
                                                                                                                        200
                                                                                                                             2011-08-18 08:5
      2011-12-09 12:50:00 UTC
                                  2011-10-05
                                                    569650
                                                                                22728
                                                                                                                             2011-10-05 12:4
  9 2011-12-09 12:50:00 UTC
                                  2011-06-24
                                                    558041
                                                                                21391
                                                                                                                             2011-06-24 13:2
  10 2011-12-09 12:50:00 UTC
                                  2011-06-20
                                                    557472
                                                                                35933
                                                                                                                             2011-06-20 13:2
  11 2011-12-09 12:50:00 UTC
                                  2011-12-09
                                                    581476
                                                                                23167
                                                                                                                             2011-12-09 08:4
  12
      2011-12-09 12:50:00 UTC
                                  2011-09-21
                                                    567526
                                                                                21524
                                                                                                                         50
                                                                                                                             2011-09-21 09:0
  13
      2011-12-09 12:50:00 UTC
                                  2011-10-26
                                                    572886
  14 2011-12-09 12:50:00 UTC
                                  2011-10-05
                                                    569653
                                                                                21376
                                                                                                                         4
                                                                                                                             2011-10-05 12:5
  15 2011-12-09 12:50:00 UTC
                                  2011-07-07
                                                    559298
                                                                                23319
                                                                                                                             2011-07-07 12:3
```

# Q. most\_recent\_date의 값 중 가장 최근 구매일은 몇 일인가요?

- 2011-12-09

이번에는 유저 별로 가장 최근에 일어난 구매 정보를 정리해 봅시다. 유저 별로 가장 큰 InvoiceDay를 찾아서 가장 최근 구매일로 저장해 주겠습니다.

```
314 SELECT
315 CustomerID,
316 MAX(DATE(InvoiceDate)) AS most_recent_purchase_date
317 FROM
318 praxis-beacon-464902-t2.modulabs_project.data2
319 GROUP BY
320 CustomerID;
```

# 쿼리 결과

작업 정도	결과	차.	E	JSON
행 //	CustomerID +	/, T	nost_r	recent_pur
1	1234	16 2	2011-0	11-18
2	1234	17 2	2011-1	2-07
3	1234	18 2	2011-0	9-25
4	1234	19 2	2011-1	1-21
5	1235	50 2	2011-0	12-02
6	1235	52 2	2011-1	1-03
7	1235	53 2	2011-0	5-19
8	1235	54 2	2011-0	14-21
9	1235	55 2	2011-0	5-09
10	1235	66 2	2011-1	1-17
11	1235	57 2	2011-1	1-06
12	1235	58 2	2011-1	2-08

다음에는 가장 최근 일자(most\_recent\_date)와 유저별 마지막 구매일(InvoiceDay)간의 차이를 계산하겠습니다.

```
322 SELECT
323 CustomerID,
324 EXTRACT(DAY FROM MAX(InvoiceDay) OVER () - InvoiceDay) AS recency
325 FROM (
326 SELECT
327 CustomerID,
328 MAX(DATE(InvoiceDate)) AS InvoiceDay
FROM praxis-beacon-464902-t2.modulabs_project.data2
330 GROUP BY CustomerID
331 );
```

# 쿼리 결과

작업 정5	결과	차트 JSON	실현
행 /,	CustomerID -	recency -	11
1	12410	30	1
2	12445	2:	2
3	12515	353	3
4	12793	334	4
5	13045	99	9
6	13075	129	9
7	13106	120	В
8	13149	14	4
9	13165	4	5

위의 코드를 더 자세히 살펴보도록 하겠습니다.

EXTRACT 함수는 날짜 또는 시간 데이터 타입에서 특정 부분을 추출하는 데 사용됩니다. 이 함수는 주로 SQL에서 날짜와 시간 연산을 수행할 때 사용되며, 날짜의 연도, 월, 일 또는 시간의 시, 분, 초 등을 추출할 수 있습니다.

쿼리문에서 EXTRACT 함수의 사용 방법을 자세히 살펴보겠습니다.

- 1. **MAX(InvoiceDay) OVER () InvoiceDay** : 각 고객(CustomerID)의 각 구매일(InvoiceDay)과 전체 데이 터셋에서의 마지막 구매일(MAX(InvoiceDay)) 간의 차이를 계산합니다.
- 2. **EXTRACT(DAY FROM ...)** : 여기에서 EXTRACT 함수는 위에서 계산된 날짜 차이에서 일(DAY) 부분만을 추출합니다. 즉 각 고객의 최근 구매일로부터 해당 구매 건의 구매일부터의 날짜 차이를 계산하는 함수 입니다.

이제 최종 데이터 셋에 필요한 데이터들을 각각 정제해서 이어붙이도록 하겠습니다. 지금까지의 결과를 user r이라는 이름의 테이블로 저장해 주세요.

```
333 CREATE OR REPLACE TABLE 'praxis-beacon-464902-t2.modulabs_project.user_r' AS
334
335 WITH last_purchase AS (
336 SELECT
337
        CustomerID,
       MAX(DATE(InvoiceDate)) AS last_invoice_day
338
339
      'praxis-beacon-464902-t2.modulabs_project.data2'
340
      GROUP BY
341
342
     CustomerID
343 ),
344
345 most_recent AS (
346 SELECT
       MAX(last_invoice_day) AS overall_last_day
347
349 | last_purchase
350 )
351
352 SELECT
     lp.CustomerID,
DATE_DIFF(mr.overall_last_day, lp.last_invoice_day, DAY) AS recency
353
354 DAT
355 FROM
356 last_purchase lp
357 CROSS JOIN
358 most_recent mr;
```

스키마	세부정보	미리트	보기
행 / C	ustomerID //	recency	//
1	15311		0
2	15694		0
3	16626		0
4	12713		0
5	15910		0
6	17315		0
7	14441		0
8	17428		0
9	12423		0
10	14422		0
11	12433		0
12	12526		0
13	17754		0
14	12985		0
15	16446		0
16	13069		0
17	16954		0
18	12662		0
19	15344		0
20	17364		0

# Frequency

Frequency를 계산하는 단계에서는 고객의 구매 빈도 또는 참여 빈도에 초점을 맞춥니다. 전체 거래 건수로 계산을 할 수도 있고, 구매한 아이템의 수량을 합하여 계산할 수도 있습니다. 예를 들어 한 명의 고객이 구매를 2번했는데 각각 아이템을 4개씩 구매한 경우, 해당 고객의 거래 건수는 2회겠지만 실제로 구매한 수량은 8개가 됩니다. 이 두가지 측면을 모두 포착하기 위해 두 개를 모두 계산해 봅시다.

# 1. 전체 거래 건수 계산

우선 각 고객의 거래 건수를 세어 봅시다. 거래 건은 InvoiceNo를 기준으로 파악하면 되기 때문에, 고객마다 고유한 InvoiceNo의 수를 세어 주겠습니다.



## 2. 구매한 아이템의 총 수량 계산

그 다음으로는 각 고객 별로 구매한 아이템의 총 수량을 더해주겠습니다.

```
395 SELECT
396 CustomerID,
397 SUM(Quantity) AS item_cnt
398 FROM
399 'praxis-beacon-464902-t2.modulabs_project.data2'
400 GROUP BY
401 CustomerID;
```

# 쿼리 결과

작업 정보	<b>로</b> 결과	차트	JSON
행 //	CustomerID *	/ item_c	ent 🕶
1	123	346	0
2	123	347	2458
3	123	148	2332
4	123	349	630
5	123	350	196
6	123	352	463
7	123	353	20
8	123	354	530
9	123	355	240
10	123	356	1573

이제 위에서 구한 '1. 전체 거래 건수 계산'과 '2. 구매한 아이템의 총 수량 계산'의 결과를 합쳐서 user\_rf라는 이름의 테이블에 저장해 주겠습니다.

```
403 CREATE OR REPLACE TABLE 'praxis-beacon-464902-t2.modulabs_project.user_rf' AS
404
405 WITH purchase_cnt AS (
     -- (1) 고객별 고유 거래 건수 계산
406
407
      SELECT
408
        CustomerID,
       COUNT(DISTINCT InvoiceNo) AS purchase_cnt
409
410
411
      'praxis-beacon-464902-t2.modulabs_project.data2'
412
     GROUP BY
413 CustomerID
414 ),
415
416 item_cnt AS (
417 -- (2) 고객별 구매 아이템 총 수량 계산
418 SELECT
419
        CustomerID,
        SUM(Quantity) AS item_cnt
420
      FROM
421
422
        'praxis-beacon-464902-t2.modulabs_project.data2'
      _praxis
423
424 CustomerID
425 )
426
427 SELECT
428 pc.CustomerID,
429 pc.purchase_cnt,
430 ic.item_cnt,
431 ur.recency
432 FROM
433 purchase_cnt AS pc
434 JOIN
     item_cnt    AS ic
ON pc.CustomerID = ic.CustomerID
435
436
437 JOIN
438 <u>praxis-beacon-464902-t2.modulabs_project.user_r</u> AS ur
439 ON pc.CustomerID = ur.CustomerID;
```

	+⊈ 공유	에서 열기 🕶	Q 쿼리 다음	er_rf	us us
E	프리뷰	테이블 탐색기	미리보기	세부정보	스키마
1	recency	item_cnt //	purchase_cnt //	CustomerID //	행 //
0		505	1	12713	1
1		79	1	14569	2
1		96	1	13298	3
1		76	1	13436	4
1		314	1	15520	5
2		256	1	15471	6
2		1404	1	15195	7
2		72	1	14204	8
3		171	1	16528	9
3		240	1	14578	10

#### Monetary

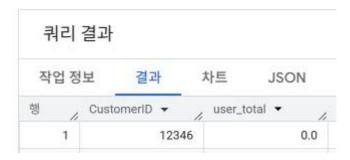
Monetary를 계산하는 단계에서는 **고객이 지불한 총 금액**에 초점을 맞춥니다. 이 때 총 지출액을 계산할수도 있고, 거래당 평균 거래 금액을 계산할 수도 있습니다. 예를 들어, 한 명의 고객이 총 2번의 구매를 했고, 그 합산 금액이 10만원인 경우, 총 지출액은 10만원, 거래당 평균 거래 금액은 5만원이 되는 것입니다.

결제한 총 금액이 높은 고객을 찾는 것도 좋지만, 한번에 많이 구매하는 고객들을 찾는 것도 굉장히 중요합니다. 결제 금액은 낮지만 구매를 자주 하는 고객과, 한번 결제할 때 큰 금액을 결제하는 고객은 분명 특성이나 행동 패턴이 다를테니까요.

이 두 가지 측면을 모두 포착하기 위해 고객별 총 지출액과 평균 거래 금액을 모두 계산해 봅시다.

## 1. 고객별 총 지출액 계산

고객별 총 지출액을 계산해 보세요. 소수점 첫째 자리에서 반올림하세요.



# 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\_cnt로 나누어서 3) user\_rfm 테이블로 저장해 봅시다.

```
449 CREATE OR REPLACE TABLE 'praxis-beacon-464902-t2.modulabs_project.user_rfm' AS
450
451
     SELECT
452
       rf.CustomerID
                                            AS CustomerID.
453
       rf.purchase_cnt,
454
       rf.item_cnt,
455
       rf.recency,
456
       ut.user_total,
457
       ROUND(ut.user_total / rf.purchase_cnt, 1) AS user_average
458
459
       `praxis-beacon-464902-t2.modulabs_project.user_rf` AS rf
460
     LEFT JOIN (
      -- (1) 고객별 총 지출액 계산
461
462
       SELECT
463
         CustomerID,
         ROUND(SUM(Quantity * UnitPrice), 1) AS user_total
464
465
         'praxis-beacon-464902-t2.modulabs_project.data'
466
467
       GROUP BY
468
        CustomerID
469
     ) AS ut
470
    ON rf.CustomerID = ut.CustomerID;
```



# RFM 통합 테이블 출력하기

위와 같은 과정을 거쳐서 생성된 최종 user\_rfm 테이블을 출력해주겠습니다.



# Q. 고유한 유저의 수는 몇 명인가요?

4,362명

# 11-8. 추가 Feature 추출

1. 구매하는 제품의 다양성

이 단계에서는 고객들의 제품 구매 행동 속 구매 제품의 다양성을 살펴보려고 합니다. 고객이 얼마나 다양한 제품들에 관심 있는 사람인지를 알게 되면, 개인 맞춤형 마케팅 전략과 추천 서비스를 계획하는 데에도 큰 도움이될 수 있습니다.

우선 1) 고객 별로 구매한 상품들의 고유한 수를 계산합니다. 높은 숫자가 나오는 것은 해당 고객이 다양한 제품들을 구매한다는 의미이며, 낮은 값이 나오는 경우 소수의 제품들만 구매한다는 것을 의미합니다.

# 이후 2) user\_rfm **테이블과 결과를 합치고**, 이를 3) user\_data**라는 이름의 테이블에 저장**하겠습니다.

```
477 CREATE OR REPLACE TABLE praxis-beacon-464902-t2.modulabs_project.user_data AS
478 WITH unique_products AS (
479
        SELECT
480
           CustomerID,
481
           COUNT(DISTINCT StockCode) AS unique_products
482
        FROM praxis-beacon-464902-t2.modulabs_project.data2
483
        GROUP BY CustomerID
484
485
     SELECT ur.*, up.* EXCEPT (CustomerID)
486 FROM praxis-beacon-464902-t2.modulabs_project.user_rfm AS ur
     JOIN unique_products AS up
488 ON ur.CustomerID = up.CustomerID;
                                                                              ■ 삭제
                     Q 쿼리
                               다음에서 열기 -
                                                         「 복사
                                                                  田 스냅샷
                                                                                        소 내보니
user data
                                              +을 공유
 스키마
           세부정보
                      미리보기
                                 테이블 탐색기 프리뷰
                                                       통계
                                                              계보
                                                                       데이터 프로필
                                                                                     데이터 품질
                                                                 user_average unique_products
                 purchase_cnt
행
       CustomerID
                               item_cnt
                                                       user_total
             17956
                            1
                                        1
                                                   249
                                                               12.8
                                                                          12.8
                                                                                          1
                                       72
    2
             17291
                            1
                                                   308
                                                              550.8
                                                                         550.8
                                                                                          1
    3
             17331
                            1
                                       16
                                                   123
                                                              175.2
                                                                          175.2
                                                                                          1
    4
             16953
                            1
                                       10
                                                    30
                                                               20.8
                                                                          20.8
                                                                                          1
             14090
                                       72
                                                   324
                                                              76.3
                                                                          76.3
    5
                            1
                                                                                          1
    6
             18068
                            1
                                        6
                                                   289
                                                              101.7
                                                                          101.7
                                                                                          1
    7
             18141
                            1
                                       -12
                                                   360
                                                              -35.4
                                                                          -35.4
                                                                                          1
                            1
                                        -1
                                                              -29.9
             16061
                                                   269
                                                                          -29.9
                                                                                          1
    9
             13366
                                                    50
                            1
                                       144
                                                               56.2
                                                                          56.2
                                                                                          1
 10
             16881
                            1
                                       600
                                                    66
                                                              432.0
                                                                         432.0
                                                                                          1
```

## 2. 평균 구매 주기

이 단계에서는 고객들의 쇼핑 패턴을 이해하는 것을 목표로 합니다. 그 중에서도 고객 별 재방문 주기를 살펴볼 것입니다. 고객들의 구매와 구매 사이의 기간이 평균적으로 몇 일인지를 보여주는 평균 일수를 계산하면, 고객이 다음 구매를 언제할지 예측하는 데에도 큰 도움이 됩니다. 평균 구매 소요 일수를 계산하고, 그 결과를 user\_data 에 통합해 줍시다.

```
CREATE OR REPLACE TABLE praxis-beacon-464902-t2.modulabs_project.user_data AS
478
     WITH purchase_intervals AS
          (2) 고객 별 구매와 구매 사이의 평균 소요 일수
479
489
       SELECT.
481
         CustomerTD
482
         CASE WHEN ROUND(AVG(interval_), 2) IS NULL THEN 0 ELSE ROUND(AVG(interval_), 2) END AS average_interval
483
       FROM (
484
            (1) 구매와 구매 사이에 소요된 일수
485
         SELECT
486
           CustomerID,
           DATE_DIFF(InvoiceDate, LAG(InvoiceDate) OVER (PARTITION BY CustomerID ORDER BY InvoiceDate), DAY) AS interval_
487
488
         FROM
489
           praxis-beacon-464902-t2.modulabs_project.data2
499
         WHERE CustomerID IS NOT NULL
491
492
      GROUP BY CustomerID
493
494
     SELECT u.*, pi.* EXCEPT (CustomerID)
     FROM praxis-beacon-464902-t2.modulabs_project.user_data AS u
     LEFT JOIN purchase_intervals AS pi
498 ON u.CustomerID = pi.CustomerID;
 user data
                         Q 쿼리
                                    다음에서 열기 マ
                                                       +일 공유
                                                                   □ 복사
                                                                              ₫ 스냅샷
                                                                                           출 삭제
                                                                                                      소 내보내기
  스키마
                          미리보기
                                                                         계보
                                                                                  데이터 프로픽
                                                                                                    데이터 포직
             세부정보
                                       테이블 탐색기 프리뷰
                                                                토계
                                                                                             unique_prod...
                                                                               user_average //
         CustomerID
                                                                 user_total
                       purchase_cnt
                                     item_cnt
                                                   recency
                                                                                                           average_interval
     1
                16953
                                  1
                                               10
                                                             30
                                                                         20.8
                                                                                        20.8
                                                                                                        1
                                                                                                                       0.0
     2
                                  1
                13270
                                              200
                                                            366
                                                                        590.0
                                                                                       590.0
                                                                                                        1
                                                                                                                       0.0
     3
                16138
                                  1
                                               -1
                                                            368
                                                                          -8.0
                                                                                        -8.0
                                                                                                        1
                                                                                                                       0.0
     4
                                                             21
                                                                        739.2
                                                                                       739.2
                                                                                                        1
                                                                                                                       0.0
                12603
                                              56
     5
                18113
                                  1
                                               72
                                                            368
                                                                         76.3
                                                                                        76.3
                                                                                                                       0.0
                                                                                                        1
                                              288
                                                             53
                                                                        417.6
                                                                                       417.6
                16737
                                                                                                                       0.0
                14090
                                  1
                                               72
                                                            324
                                                                         76.3
                                                                                        76.3
                                                                                                        1
                                                                                                                       0.0
     8
                14119
                                  1
                                               -2
                                                            354
                                                                         -19.9
                                                                                       -19.9
                                                                                                        1
                                                                                                                       0.0
                                  1
     9
                13120
                                               12
                                                            238
                                                                         30.6
                                                                                        30.6
                                                                                                        1
                                                                                                                       0.0
     10
                13135
                                             4300
                                                            196
                                                                       3096.0
                                                                                      3096.0
                                                                                                        1
                                                                                                                       0.0
```

# 3. 구매 취소 경향성

이 단계에서는 **고객의 취소 패턴**을 더 깊이 있게 파고 들어 고객을 세그먼테이션할 때 사용하려고 합니다. 아래와 같은 특징들을 추가해 보겠습니다.

#### 1. 취소 빈도(cancel\_frequency)

취소 빈도는 고객 별로 취소한 거래의 총 횟수입니다.

취소 빈도를 이해하면 거래를 취소할 가능성이 높은 고객을 식별할 수 있습니다. 취소 빈도는 불만족의 정도이거나 다른 문제에 대한 지표일 수 있습니다. 따라서 취소 빈도를 이해함으로써 거래 취소 횟수를 줄이고 고객 만족도를 높이기 위한 전략을 세울 수 있습니다.

#### 2. 취소 비율(cancel rate)

취소 비율은 각 고객이 한 모든 거래 중에서 취소를 한 거래의 비율입니다.

취소 비율은 특정 고객이 원래 취소를 잘 하는 고객인지와 같은 고객의 특징을 잡아내기 위한 지표입니다. 이런 특성을 식별함으로써 고객의 쇼핑 경험을 개선하고 취소 비율을 감소시키기 위해 어떤 고객 대상군을 공략해야 할지에 대한 실마리를 얻을 수 있습니다. 취소 빈도와 취소 비율을 계산하고 그 결과를 user\_data에 통합해 줍시다. 취소 비율은 소수점 두번째 자리까지 구해 주세요.

478 479 V			2000				0.00	100	data' AS
179 1	HOUSE TO THE REAL PROPERTY OF THE PERTY OF T								
	WITH Trans	sactionIr	nto AS (						
480	SELECT	70							
481	Custor					10	. 750		
482			InvoiceN					nsaction	
483		IF (STARTS	S_WITH(Inve	oiceNo, C	( ) )	AS (	cance1_f	requency	
484	FROM		45 4000 +4			1			
485	GROUP BY		n-464902-t2	2.modulabs	_project.	data			
186	100								
487	Custor	nerio							
488 ) 489	)								
	SELECT								
	rf.*,								
	t.total	trancaci	tione						
	t.cance.								
194			frequency .	/ t total	transacti	one * 100	2) 45	cancel r	ate
	FROM	.cancer_	requeitey	c.cocur_	er ansacti	100	2/ 10	cancer_n	acc
196		-heacon-4	464902-t2.r	modulahs m	roject us	er rfm` As	S rf		
	LEFT JOIN	beacon	TO TOOL CE	"oddrabo_p	n ojece.ac				
	EFF JULN					1750			
		tionInfo	AS t						
198		tionInfo	AS t			1980			
198	Transact ON	176-1	120	rID;		- 100 C			
98 199 (	Transact ON	omerID =	AS t t.Custome 다음에서 열기 ㅜ	100	「□ 복사	스냅샷 🍵 삭	제 소 내5	보내기	
98 199 (	Transact ON rf.Custo	omerID =	t.Custome	100		스냅샷 <b>흡</b> 삭	제 소 내5	ication i	
198 199 ( 190 <b>1</b> use 스키마	Transaction Transa	omerID = Q 쿼리	t.Custome 다음에서 열기 ▼ 테이블 탐색기	+ <b>2</b> 공유 프리뷰 통7	계 계보	데이터 프로필	데이터 품	ication i	cancel_rate
198 199 ( 500 <b>I</b> use	Transaction Transa	omerID = Q 쿼리 미리보기	t.Custome 다음에서 열기 ▼ 테이블 탐색기	+ <b>2</b> 공유 프리뷰 통7	계 계보	데이터 프로필	데이터 품	질	
198 199 ( 500   로 use	Transaction  on  rf.Custo er_data  MPSE  CustomerID	omerID = Q 쿼리 미리보기	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt	* 공유 프리뷰 통기 ecency / us	례 계보 ser_total us	데이터 프로필 ser_average / tota	데이터 품 al_transac	질 cancel_frequ	cancel_rate 0.0
198 199 ( 500 use 스키마	Transaction  rf.Custo er_data 세부정보  CustomerID / 1 13844	omerID = Q 쿼리 미리보기 ourchase_cnt //	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt / r	# 공유 프리뷰 통2 ecency ##	비 계보 ser_total us 361.6	데이터 프로필 ser_average // tota 361.6	데이터 품 al_transac	질 cancel_frequ/	0.0
198 199 (500 ■ use △키마	Transaction  rf.Custo er_data 세부정보  CustomerID # 13844 # 15043	omerID = Q 쿼리 미리보기 ourchase_cnt // 1	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt r 196 300	** 39	계 계보 ser_total us 361.6 537.0	데이터 프로필 ser_average / tota 361.6 537.0	데이터 품 al_transac	질 cancel_frequ 0 0	0.0
98   99 (600   use 스키마	Transaction  rf.Custo er_data  MPSE  CustomerID  13844  15043  16318	omerID = Q 쿼리 미리보기 ourchase_cnt // 1	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt	** 39	에 계보 ser_total us 361.6 537.0 328.1	데이터 프로필 ser_average tota 361.6 537.0 328.1	데이터 품 al_transac,	질 cancel_frequ 0 0	0.1 0.1 0.1
198 199 ( 199 ( 1	Transaction  rf.Custo er_data 세부정보  CustomerID	omerID = Q 쿼리 미리보기 ourchase_cnt // 1 1	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt / 196 300 199 362 82	*** 38	期 계보 361.6 537.0 328.1 754.6 243.9	데이터 프로필 361.6 537.0 328.1 754.6 243.9	데이터 품 1 1 1 1 1	질 cancel_frequ 0 0 0	0.0 0.1 0.1 0.1
98   99 (600   use	Transaction  rf.Custo er_data  All PSE  13844  15043  16318  18104  13103  15385	omerID = Q 쿼리 미리보기 purchase_cnt // 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt r 196 300 199 362 82 161	** 39 ** 37 39 40	利보 ser_total us 361.6 537.0 328.1 754.6 243.9 316.9	데이터 프로필 361.6 537.0 328.1 754.6 243.9 316.9	데이터 품 al_transac/ 1 1 1 1 1	질 cancel_frequ / 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	0.4 0.4 0.4 0.4 0.4
98   99 (600   use 스키마 1 2 3 4 5 6 7 7	Transaction  rf.Custo er_data 세부정보  CustomerID 13844 15043 16318 18104 13103 15385 13845	omerID = Q 쿼리 미리보기 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt 196 300 199 362 82 161 182	** 39 40 64	制 계보 ser_total us 361.6 537.0 328.1 754.6 243.9 316.9 313.0	데이터 프로필 361.6 537.0 328.1 754.6 243.9 316.9	데이터 품 al_transac/ 1 1 1 1 1 1	짇 O O O O O O O O O O O O O O O O O O O	0.1 0.3 0.1 0.1 0.1
98   99   0   99   0   99   0   99   0   99   0   99   0   99   0   99   0   99	Transaction  rf.Custo er_data  // MP정보  CustomerID / F  13844  15043  16318  18104  13103  15385  13845  13514	omerID = Q 쿼리 미리보기	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_ont 196 300 199 362 82 161 182 40	** 공유  프리뷰 통기  eccency us  11  31  35  37  39  40  64  73	制 계보 361.6 537.0 328.1 754.6 243.9 316.9 313.0	데이터 프로필 361.6 537.0 328.1 754.6 243.9 316.9 313.0	데이터 품 al_transac	질 cancel_frequ / 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	0.1 0.1 0.1 0.1 0.1 0.1
198   199 (199   199 (199 (199 (199 (199 (1	Transaction  rf.Custo er_data  All P B L  13844  15043  16318  18104  13103  15385  13845  13514  15325	omerID = Q 쿼리 미리보기 ourchase_cnt // 1 1 1 1 1 1 1 1	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_cnt / 196 300 199 362 82 161 182 40 231	** 38	明 계보 361.6 537.0 328.1 754.6 243.9 316.9 313.0 152.2 162.3	데이터 프로필 361.6 537.0 328.1 754.6 243.9 316.9 313.0 152.2 162.3	데이터 품 1 1 1 1 1 1 1 1 1	Dancel_frequ	0.1 0.1 0.1 0.1 0.1 0.1
198   199 (199   199 (199 (199 (199 (199 (1	Transaction  rf.Custo er_data  // MP정보  CustomerID / F  13844  15043  16318  18104  13103  15385  13845  13514	omerID = Q 쿼리 미리보기	t.Custome 다음에서 열기 ▼ 테이블 탐색기 item_ont 196 300 199 362 82 161 182 40	** 공유  프리뷰 통기  eccency us  11  31  35  37  39  40  64  73	制 계보 361.6 537.0 328.1 754.6 243.9 316.9 313.0	데이터 프로필 361.6 537.0 328.1 754.6 243.9 316.9 313.0	데이터 품 al_transac	질 cancel_frequ / 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	0.0 0.1 0.1 0.1 0.1 0.1 0.1 0.1

지금까지 추가적인 feature를 추출하기 위한 작업을 진행하였습니다.

다양한 컬럼들을 활용하여 고객의 구매 패턴과 선호도를 보다 심층적으로 이해할 수 있도록 최종적으로 user\_data를 출력해보겠습니다.

```
477 SELECT
478 *
479 FROM
480 | `praxis-beacon-464902-t2.modulabs_project.user_data`;
```

작업 정	보 결과	차트 JSON	실행 세부정보	실행 그래프					
행 //	CustomerID ▼	purchase_cnt • //	item_cnt ▼	recency •	user_total ▼ //	user_average ▼ //	total_transactions	cancel_frequency •	cancel_rate -
1	13844	1	196	11	361.6	361.6	1	0	0.
2	15043	1	300	31	537.0	537.0	1	0	0.
3	16318	1	199	35	328.1	328.1	1	0	0.
4	18104	1	362	37	754.6	754.6	1	0	0.
5	13103	1	82	39	243.9	243.9	1	0	0.
6	15385	1	161	40	316.9	316.9	1	0	0.
7	13845	1	182	64	313.0	313.0	1	0	0