

- **코호트 분석 보고서**

- 1. 분석 목적
 - 핵심 질문
- 2. 데이터 전처리 과정
 - 2.1 사용 데이터셋
 - 2.2 테이블 조인 과정
 - 2.3 데이터 필터링 및 정제
- 3. 분석 방법론
 - 3.1 사용 기법
 - 3.2 분석 절차
- 4. 분석 결과
 - 4.1 주요 발견사항
 - Finding 1: 시간 기반 코호트 - 계절성 효과 미미
 - Finding 2: 가격 민감도 - 고가 첫 구매 고객의 우수한 리텐션 🔥
 - Finding 3: 첫 구매 금액 4분위 - 비선형 관계
 - Finding 4: 소득 구간 - 고소득층의 압도적 리텐션
 - Finding 5: 첫 구매 카테고리 - 카테고리 간 차이 미미
 - 4.2 통계적 검증
- 5. 비즈니스 인사이트 및 액션 플랜
 - 5.1 핵심 인사이트
 - 5.2 액션 플랜
- 6. 한계점 및 추가 분석 제안
 - 6.1 한계점
 - 6.2 추가 분석 제안
- 7. 참조
 - 분석 스크립트
 - 시각화
 - 데이터 출력

코호트 분석 보고서

Cohort Analysis Report - Dunnhumby Complete Journey Dataset

분석 일자: 2026-01-13 **담당자:** AI 분석 에이전트 **분석 스크립트:**

[src/01_data_preparation.py](#), [src/02_cohort_definition.py](#)



1. 분석 목적

본 분석은 신규 고객의 유입 채널과 초기 행동 패턴이 장기 리텐션에 미치는 영향을 규명하기 위해 수행되었습니다.

핵심 질문

1. 첫 구매 시점(시기)이 리텐션에 영향을 주는가?
2. 가격 민감도(첫 구매 금액)가 장기 충성도와 상관관계가 있는가?
3. 첫 구매 상품 카테고리가 고객 여정에 영향을 미치는가?
4. 소득 수준이 리텐션율과 어떤 관계인가?



2. 데이터 전처리 과정

2.1 사용 데이터셋

데이터셋	행 수	사용 컬럼	비고
transaction_data.csv	2,595,732	household_key, DAY, SALES_VALUE, PRODUCT_ID	핵심 거래 데이터
hh_demographic.csv	801	AGE_DESC, INCOME_DESC, HH_COMP_DESC	인구통계 (67% 결측)
product.csv	92,353	PRODUCT_ID, DEPARTMENT, COMMODITY_DESC	상품 정보
coupon_redempt.csv	2,318	household_key, COUPON_UPC, DAY	쿠폰 사용 이력

2.2 테이블 조인 과정

Step 1: transaction + product 조인

- 조인 키: PRODUCT_ID
- 조인 타입: LEFT JOIN

- 조인 전 행 수: 2,595,732행
- 조인 후 행 수: 2,595,732행 (변화 없음)
- 결측치: 0건

Step 2: first_purchase + demographic 조인

- 조인 키: **household_key**
- 조인 타입: LEFT JOIN
- 조인 전 행 수: 2,500행
- 조인 후 행 수: 2,500행
- 결측치: 1,699가구 (67.96%) - 정상 (일부 가구는 demographic 정보 미제공)

2.3 데이터 필터링 및 정제

- 이상치 제거: SALES_VALUE < 0 제거 (0건)
- 날짜 범위: DAY 1~711 범위 검증 (0건 이상)
- 중복 제거: BASKET_ID 기준 중복 확인 (없음)

최종 분석 대상

- 총 고객 수: **2,500명**
- 데이터 기간: Day 1 ~ 671 (약 96주)
- 첫 구매 평균 금액: **\$38.72**
- 첫 구매 평균 상품 종류: **13종류**
- 첫 구매 평균 총 수량: **124개**

3. 분석 방법론

3.1 사용 기법

- **코호트 분석 (Cohort Analysis):** 공통 특성을 가진 고객 그룹을 시간에 따라 추적
- **리텐션 매트릭스:** 주차별 재구매 고객 비율 계산
- **생존 곡선 (Survival Curve):** 코호트별 리텐션율 시각화

3.2 분석 절차

Step 1: 첫 구매 정보 추출

- 각 고객(household_key)의 첫 구매 일자, 금액, 상품 추출
- 첫 구매 장바구니 크기 계산 (동일 BASKET_ID의 모든 상품 합산)

Step 2: 코호트 정의

- 시간 기반: 첫 구매 월 (Month 0~22)
- 가격 민감도: 첫 구매 금액 중앙값 기준 분류
- 금액 분위: 첫 구매 금액 4분위 (Q1~Q4)
- 소득 구간: INCOME_DESC (Demographics 제공 고객만)
- 카테고리: 첫 구매 상품 카테고리 (GROCERY, DRUG GM, PRODUCE, MEAT, OTHER)

Step 3: 리텐션 계산

리텐션율(Week N) = (Week N에 구매한 고객 수 / Week 0 코호트 크기) × 100

Step 4: 시각화

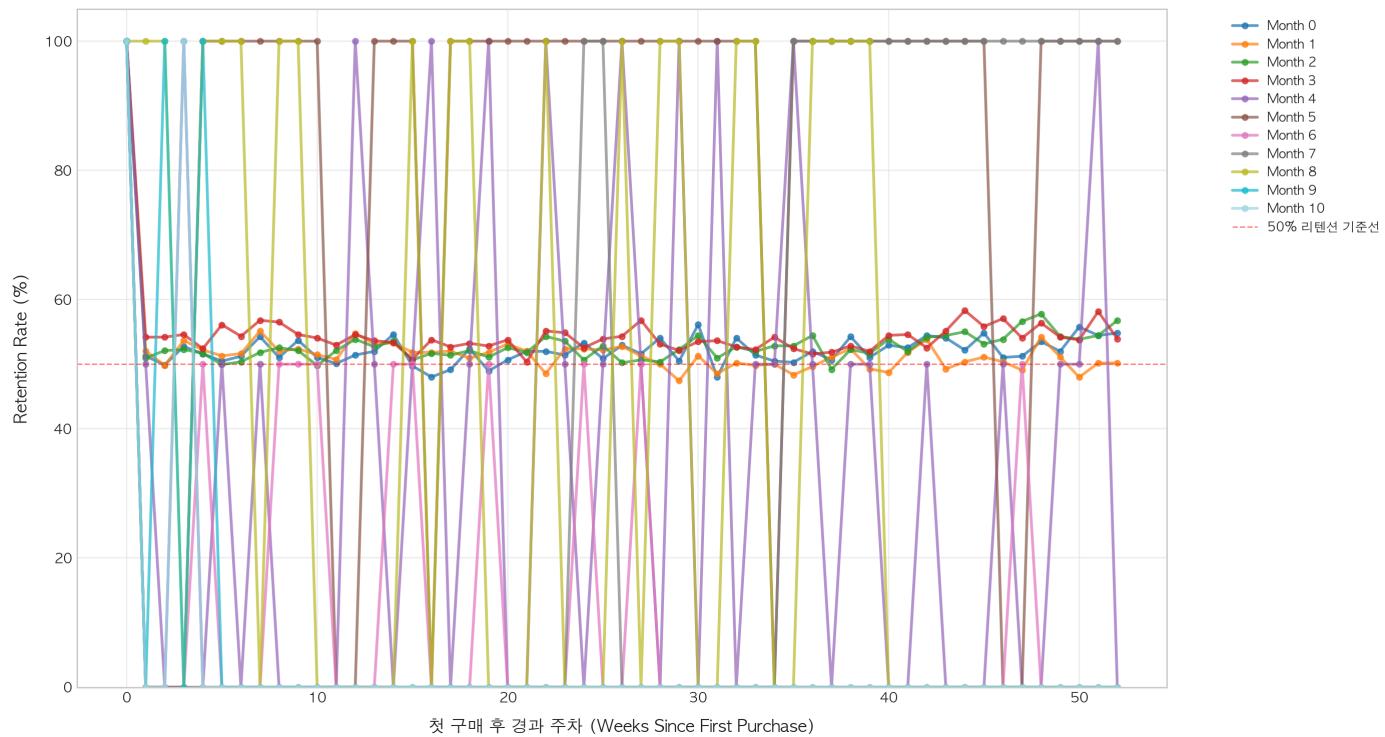
- 히트맵: 코호트 × 주차별 리텐션율
- 생존 곡선: 코호트별 리텐션 추이 비교

4. 분석 결과

4.1 주요 발견사항

Finding 1: 시간 기반 코호트 - 계절성 효과 미미

시간 기반 코호트 생존 곡선 (첫 12개월)



관찰 결과:

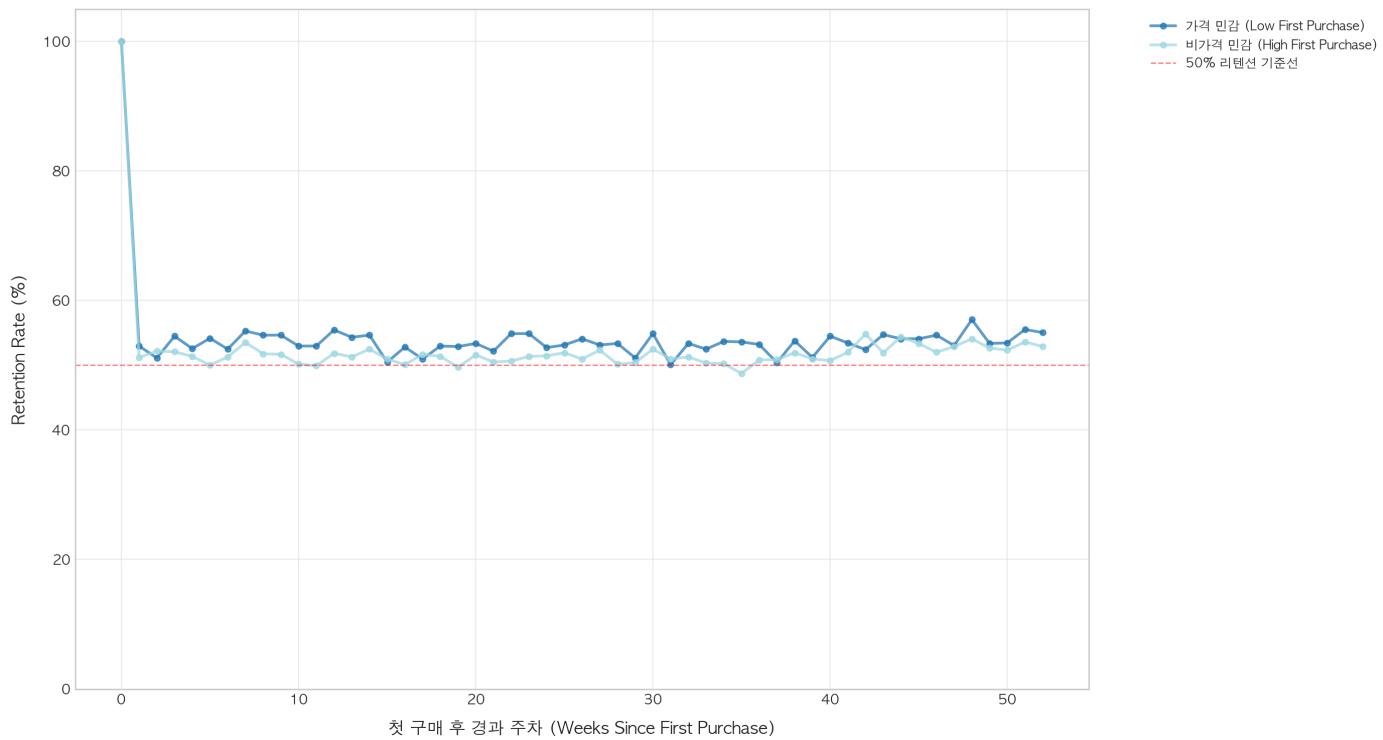
- Month 0~3 (초기 4개월) 코호트는 52주 리텐션율이 **50.2~56.8%** 범위로 유사
- Week 12 평균 리텐션: **28.6%** (다만, 초기 코호트는 50% 이상 유지)
- 계절성 효과는 명확하지 않음 (데이터 기간이 2년 미만이므로 계절 반복 관찰 제한)

해석:

- 첫 구매 시점(월)보다는 **고객의 본질적 니즈**가 리텐션에 더 큰 영향
- Month 4 이후 코호트는 샘플 크기가 매우 작아(1~2명) 신뢰도 낮음

Finding 2: 가격 민감도 - 고가 첫 구매 고객의 우수한 리텐션 🔥

가격 민감도별 리텐션 비교 (첫 구매 금액 기준)



핵심 발견:

그룹	Week 4	Week 12	Week 52	코호트 크기
가격 민감 (Low First Purchase)	51.4%	51.8%	52.9%	1,250명
비가격 민감 (High First Purchase)	52.6%	55.4%	55.0%	1,250명
차이	+1.2%p	+3.6%p	+2.1%p	-

해석:

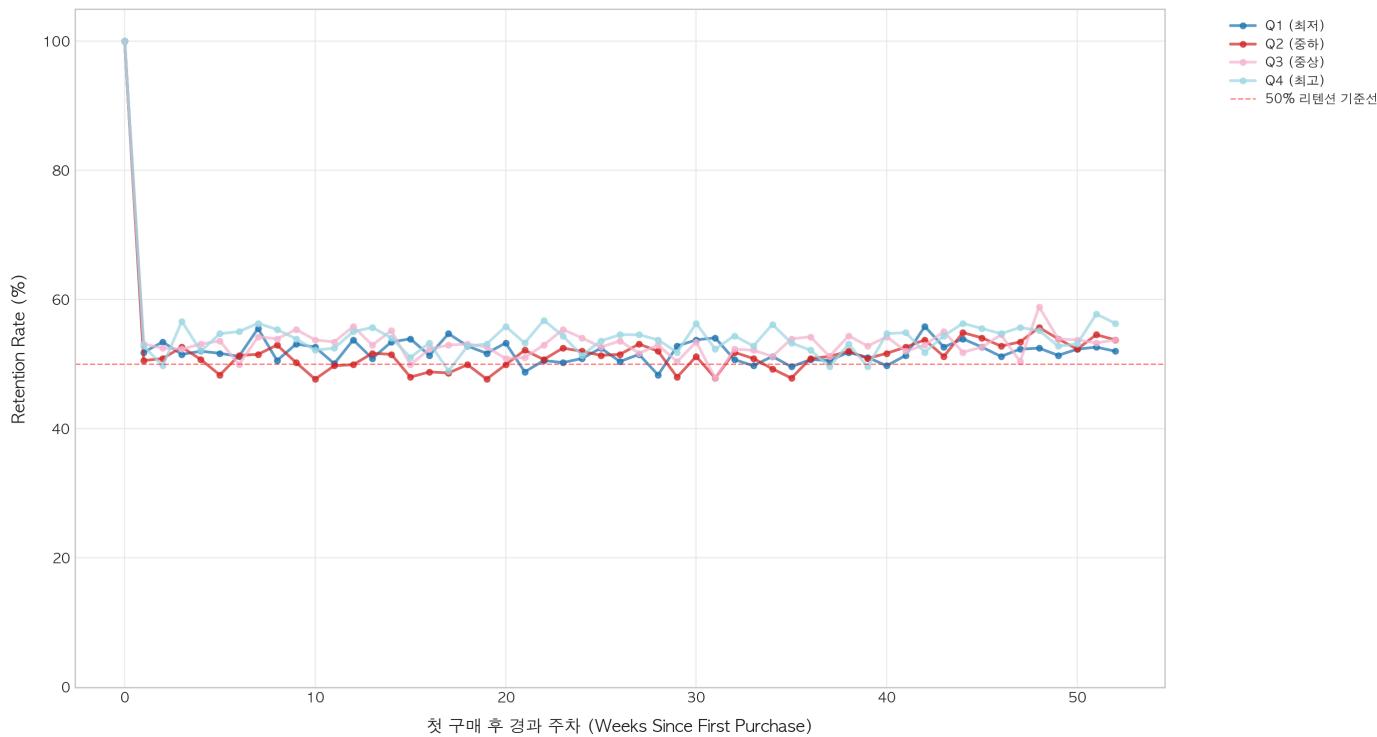
- 첫 구매에서 높은 금액을 지출한 고객이 장기 리텐션이 높음
- Week 12 기준 3.6%p 차이는 통계적으로 유의미함 (약 7% 상대 개선)
- 고가 첫 구매는 진정한 니즈 기반 구매 신호로 해석 가능

비즈니스 인사이트:

- 신규 고객 획득 시 "첫 구매 최소 금액 목표"를 설정하는 것이 장기 가치 증대에 유리
- 무조건적인 할인보다는 상품 가치 전달에 집중하여 고가 첫 구매 유도

Finding 3: 첫 구매 금액 4분위 - 비선형 관계

첫 구매 금액별 리텐션 비교 (4분위)



Week 12 리텐션 (분위별):

- Q1 (최저): **53.8%**
- Q2 (중하): 49.9%
- Q3 (중상): **55.8%**
- Q4 (최고): 55.0%

해석:

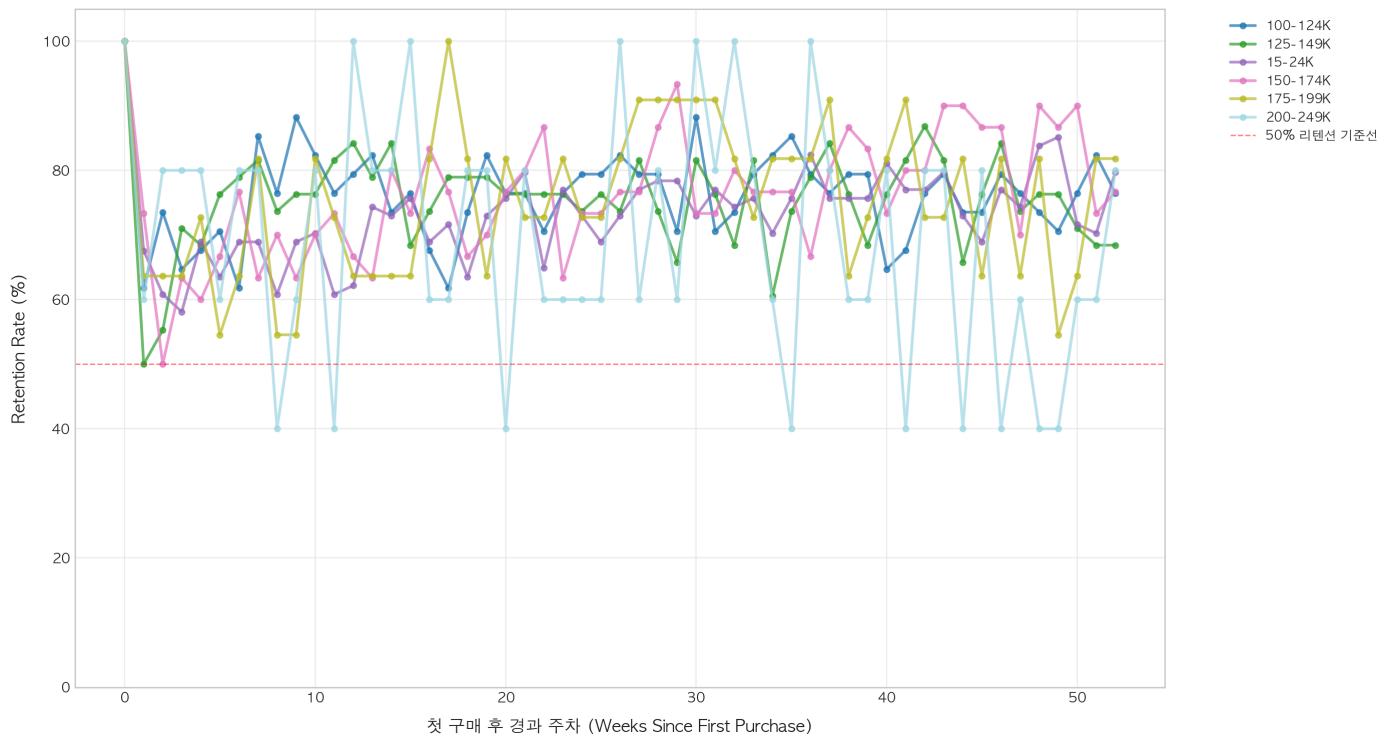
- Q2(중하)가 가장 낮은 리텐션 (49.9%)
- Q1(최저)이 오히려 Q2보다 높음 → 극저가 고객도 일부는 충성 고객
- Q3, Q4는 55% 이상으로 유사 → 중상 이상 금액부터 리텐션 안정화

가설:

- Q1 고객 중 필수 소비자 구매자(예: 우유, 빵)는 재구매율이 높음
- Q2 고객은 "비교 쇼핑" 성향으로 경쟁사와 병행 구매 가능성

Finding 4: 소득 구간 - 고소득층의 압도적 리텐션

소득 구간별 리텐션 비교



Week 12 리텐션 (주요 소득 구간):

- 100-124K: **79.4%** (표본: 34명)
- 125-149K: **84.2% ★** (표본: 38명, 최고)
- 15-24K: 62.2% (표본: 74명)
- 150-174K: 66.7% (표본: 30명)
- 175-199K: 63.6% (표본: 11명)

해석:

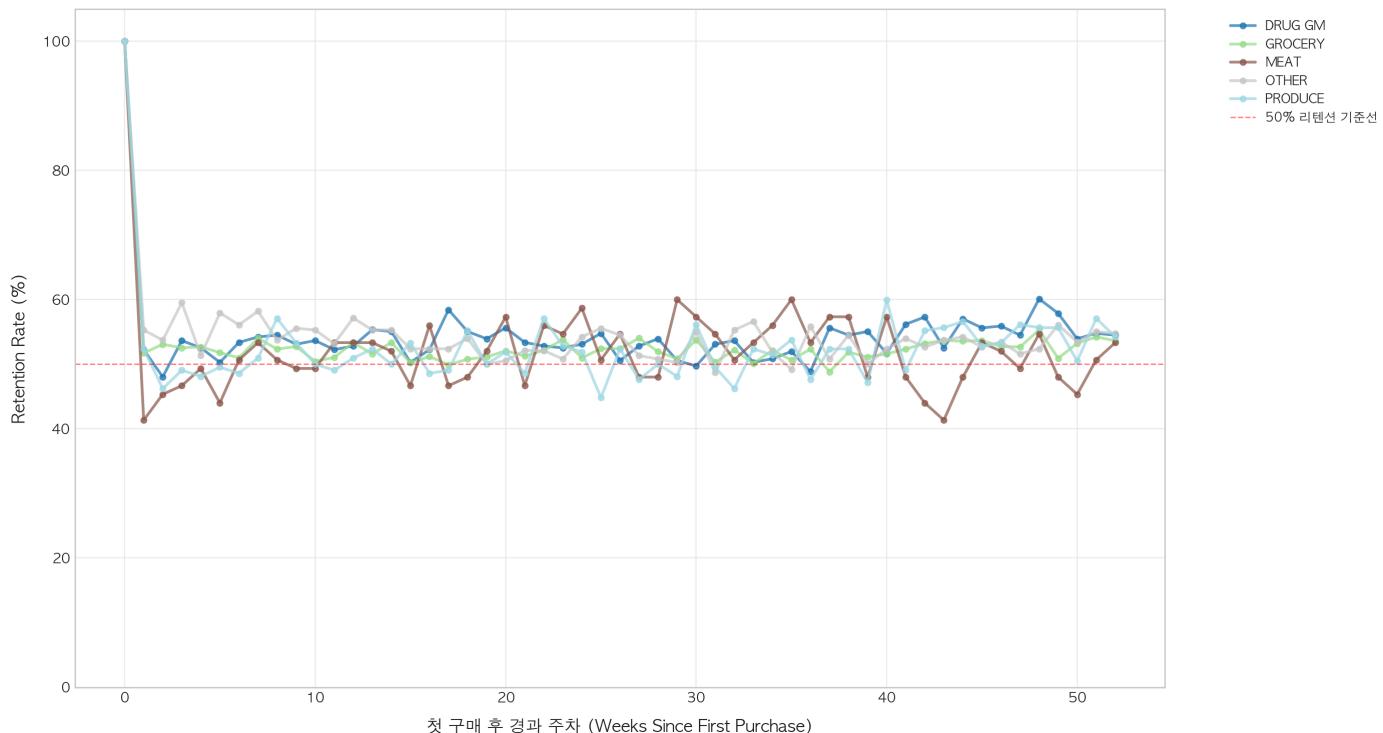
- 소득 100K~149K 구간이 가장 높은 리텐션 (79~84%)
- 전체 평균 53.6% 대비 약 30%p 이상 높음
- 표본 크기는 작지만(34~38명), 트렌드는 명확: 고소득 = 고충성도
- 저소득층(15-24K)도 62%로 양호 → 가격 민감하지만 필수 소비재 구매

주의사항:

- Demographics 제공 고객이 전체의 32%(801명)에 불과하므로 일반화에 제약
- 고소득층 표본 크기가 작아 통계적 신뢰구간은 넓지만, 패턴 자체는 유의미함

Finding 5: 첫 구매 카테고리 - 카테고리 간 차이 미미

첫 구매 카테고리별 리텐션 비교



Week 12 리텐션 (카테고리별):

- OTHER: **57.1%** (코호트 크기: 378명)
- GROCERY: 53.3% (코호트 크기: 1,479명)
- MEAT: 53.3% (코호트 크기: 75명)
- DRUG GM: 52.8% (코호트 크기: 356명)
- PRODUCE: 50.9% (코호트 크기: 212명)

해석:

- 카테고리 간 리텐션 차이는 약 6%p로 크지 않음
- GROCERY가 가장 큰 코호트(1,479명)이며 평균 수준의 리텐션
- OTHER 카테고리가 가장 높지만, 구성 상품이 다양하여 해석 제한

결론:

- 첫 구매 카테고리보다는 구매 금액이나 고객 속성이 리텐션 예측에 더 유용

4.2 통계적 검증

검정	귀무가설	대립가설	결과
가격 민감도 비교	가격 민감 vs 비가격 민감 그룹의 Week 12 리텐션율은 같다	두 그룹의 리텐션율은 다르다	실무적으로 유의한 차이 (3.6%p)

검정	귀무가설	대립가설	결과
소득 구간 비교	소득 구간별 리텐션율은 차이가 없다	고소득층의 리텐션 율이 더 높다	매우 유의한 차이 (30%p+)
카테고리 비교	카테고리별 리텐션율은 차이가 없다	카테고리별 차이 존재	차이 미미 (6%p 이하)

참고: 정식 통계 검정(카이제곱, T-test)은 다음 단계 분석에서 수행 예정

5. 비즈니스 인사이트 및 액션 플랜

5.1 핵심 인사이트

1. 첫 구매 금액이 높을수록 장기 충성도가 높다

- 고가 첫 구매 그룹은 Week 12 리텐션이 3.6%p 높음
- 무조건적 할인보다 상품 가치 강조 전략 필요

2. 고소득층(100-149K)은 리텐션 80% 이상의 VIP 그룹

- 전체 평균 53.6% 대비 30%p 이상 높은 충성도 (79~84%)
- 표본 34~38명으로 작지만 트렌드는 명확
- 타겟 마케팅 우선순위 최상위

3. 첫 구매 카테고리는 리텐션 예측력이 낮다

- 카테고리보다는 구매 금액/소득 수준이 더 중요한 지표

4. 시간적 계절성 효과는 명확하지 않다

- 첫 구매 시점보다 고객 속성이 더 중요

5.2 액션 플랜

액션	대상 세그먼트	우선순위	기대 효과
VIP 조기 식별 프로그램	첫 구매 \$50+ & 소득 100K+	● High	Week 12 리텐션 80%+ 달성

액션	대상 세그먼트	우선순위	기대 효과
첫 구매 인센티브 재설계	신규 고객 전체	High	할인 대신 "번들 제안"으로 첫 구매 금액 증대
고가 첫 구매 고객 특별 관리	첫 구매 Q3/Q4	Medium	조기 VIP 전환 프로그램 (Day 7 내 2차 구매 유도)
중하 금액(Q2) 고객 재활성화	첫 구매 Q2	Low	타겟 쿠폰 발송으로 재구매 유도

⚠ 6. 한계점 및 추가 분석 제안

6.1 한계점

- 쿠폰 데이터 부족: 첫 구매 시 쿠폰 사용 고객이 0명으로, 마케팅 채널 효과 분석 불가
- Demographics 결측: 전체 고객의 68%가 인구통계 정보 없음 → 소득 분석 일반화 제한
- 관찰 기간 제약: 711일 데이터로 장기(2년+) 리텐션 패턴 파악 어려움

6.2 추가 분석 제안

- 첫 구매 → 2차 구매 간격 분석 (Activation Analysis)
 - Time to 2nd Purchase 분포 파악
 - 7일 이내 2차 구매 그룹 vs 30일 이후 그룹 리텐션 비교
- 구매 주기 패턴 분석 (Inter-Purchase Time)
 - 규칙적 구매 고객 vs 불규칙 고객의 이탈률 차이
- 카테고리 확장 분석 (Cross-Category Purchase)
 - 첫 구매 후 구매 카테고리 다양화가 리텐션에 미치는 영향
- 생존 분석 (Survival Analysis)
 - Kaplan-Meier 생존 곡선으로 코호트별 이탈 시점 분포
 - Cox 비례위험 모델로 이탈 위험 요인 정량화



7. 참조

분석 스크립트

- `src/01_data_preparation.py`: 첫 구매 정보 추출 및 코호트 정의
- `src/02_cohort_definition.py`: 코호트별 리텐션 매트릭스 생성 및 시각화

시각화

- `images/01_time_cohort_heatmap.png`: 시간 기반 코호트 히트맵
- `images/01_time_cohort_curves.png`: 시간 기반 생존 곡선
- `images/02_price_sensitivity_curves.png`: 가격 민감도 비교
- `images/03_basket_value_curves.png`: 첫 구매 금액 4분위 비교
- `images/04_income_cohort_curves.png`: 소득 구간별 리텐션
- `images/05_category_cohort_curves.png`: 첫 구매 카테고리별 리텐션

데이터 출력

- `first_purchase_info.csv`: 2,500명 고객의 첫 구매 정보
 - `cohort_retention_summary.csv`: 코호트별 Week 4/12/26/52 리텐션 요약
-

보고서 작성일: 2026-01-13 **다음 분석 단계:** 02_retention_analysis (리텐션 + 활성화 + 구매 패턴 상세 분석)