

```

In [12]: import pandas as pd
import numpy as np
from scipy.stats import ttest_ind

# CSV 불러오기
df = pd.read_csv("C:/Users/HomePc/Desktop/review.csv")
df.dropna(subset=['review_text', 'gender', 'review_date', 'score', 'year'], inplace=True)
df['review_date'] = pd.to_datetime(df['review_date'], errors='coerce')
df.dropna(subset=['review_date'], inplace=True)

# -----
# [검정 a] 성별에 따른 영화 평가 차이
# -----

print("★ [검정 a] 성별에 따른 영화 평가 차이")
print("[제반조건]")
print("- 두 집단(남성과 여성)은 서로 독립적이다.")
print("- 표본 수는 각각 10명으로 작다.")
print("- 등분산성 여부를 알 수 없으므로 등분산을 가정하지 않는다.")
print("→ 따라서, Welch의 t-검정(독립 2표본, 비모수적 등분산 가정 없음)을 사용함.\n")

male_reviews = df[df['gender'] == 'M'].sample(10, random_state=1)
female_reviews = df[df['gender'] == 'F'].sample(10, random_state=1)

stat_a, pval_a = ttest_ind(male_reviews['score'], female_reviews['score'], equal_var=False)

print("귀무가설(H0): 남성과 여성의 평균 평점에는 차이가 없다.")
print("대립가설(H1): 남성과 여성의 평균 평점에는 차이가 있다.")
print(f"→ 검정통계량: {stat_a:.3f}, p값: {pval_a:.4f}")

if pval_a < 0.05:
    print("→ p값이 0.05 미만이므로 귀무가설을 기각하고, 대립가설을 채택한다.")
    print("→ 성별에 따라 영화 평가에 유의미한 차이가 있다고 본다.\n")
else:
    print("→ p값이 0.05 이상이므로 귀무가설을 채택한다.")
    print("→ 성별에 따른 영화 평점에 유의미한 차이는 없다고 본다.\n")

# -----
# [검정 b] 부정 리뷰의 연도별 평가 차이 (2024 vs 2025)
# -----

print("★ [검정 b] 부정 리뷰의 2024년 vs 2025년 평가 비교")
print("[제반조건]")
print("- 두 집단(2024년과 2025년)은 서로 독립적이다.")
print("- 각 집단에서 100개의 부정 리뷰를 임의 추출한다.")
print("- 등분산을 가정할 수 없으므로 Welch의 t-검정을 사용함.\n")

neg_reviews = df[df['sentiment'] == 'Negative']
neg_sample = neg_reviews.sample(100, random_state=2)

neg_2024 = neg_sample[neg_sample['year'] == 2024]['score']
neg_2025 = neg_sample[neg_sample['year'] == 2025]['score']

stat_b, pval_b = ttest_ind(neg_2024, neg_2025, equal_var=False, nan_policy='omit')

print("귀무가설(H0): 2024년과 2025년 부정 리뷰의 평균 평점에는 차이가 없다.")
print("대립가설(H1): 2024년과 2025년 부정 리뷰의 평균 평점에는 차이가 있다.")
print(f"→ 검정통계량: {stat_b:.3f}, p값: {pval_b:.4f}")

```

```

if pval_b < 0.05:
    print("⇒ p값이 0.05 미만이므로 귀무가설을 기각하고 대립가설을 채택한다.")
    print("⇒ 2024년과 2025년의 부정 평가에 유의미한 차이가 있다고 본다.\n")
else:
    print("⇒ p값이 0.05 이상이므로 귀무가설을 채택한다.")
    print("⇒ 두 연도 간 부정 리뷰의 평가에는 유의미한 차이가 없다고 본다.\n")

# -----
# [검정 c] 긍정 리뷰에서 표현("좋다" vs "그저 그렇다")에 따른 평가 차이
# -----

print("★ [검정 c] 긍정 리뷰 표현에 따른 평점 차이")
print("[제반조건]")
print("- 두 집단은 서로 독립적이다.")
print("- 표본 수는 각각 15개이며, 등분산을 가정할 수 없다.")
print("⇒ 따라서 Welch의 t-검정을 수행함.\n")

positive_reviews = df[df['sentiment'] == 'Positive']

group_good = positive_reviews[positive_reviews['review_text'].str.contains("좋다")]
group_soso = positive_reviews[positive_reviews['review_text'].str.contains("그저

stat_c, pval_c = ttest_ind(group_good['score'], group_soso['score'], equal_var=F

print("귀무가설(H0): '좋다'와 '그저 그렇다'를 포함한 긍정 리뷰의 평균 평점에는 차
print("대립가설(H1): 두 표현을 포함한 긍정 리뷰의 평균 평점에는 차이가 있다.")
print(f"⇒ 검정통계량: {stat_c:.3f}, p값: {pval_c:.4f}")

if pval_c < 0.05:
    print("⇒ p값이 0.05 미만이므로 귀무가설을 기각하고 대립가설을 채택한다.")
    print("⇒ 표현 방식에 따라 긍정 리뷰의 평가에 유의미한 차이가 있다고 본다.\n")
else:
    print("⇒ p값이 0.05 이상이므로 귀무가설을 채택한다.")
    print("⇒ 표현 방식에 따른 긍정 평가에는 차이가 없다고 본다.\n")

```

🔴 [검정 a] 성별에 따른 영화 평가 차이

[제반조건]

- 두 집단(남성과 여성)은 서로 독립적이다.
  - 표본 수는 각각 10명으로 작다.
  - 등분산성 여부를 알 수 없으므로 등분산을 가정하지 않는다.
- 따라서, Welch의 t-검정(독립 2표본, 비모수적 등분산 가정 없음)을 사용함.

귀무가설( $H_0$ ): 남성과 여성의 평균 평점에는 차이가 없다.

대립가설( $H_1$ ): 남성과 여성의 평균 평점에는 차이가 있다.

→ 검정통계량: -1.791, p값: 0.0907

⇒ p값이 0.05 이상이므로 귀무가설을 채택한다.

⇒ 성별에 따른 영화 평점에 유의미한 차이는 없다고 본다.

🔴 [검정 b] 부정 리뷰의 2024년 vs 2025년 평가 비교

[제반조건]

- 두 집단(2024년과 2025년)은 서로 독립적이다.
- 각 집단에서 100개의 부정 리뷰를 임의 추출한다.
- 등분산을 가정할 수 없으므로 Welch의 t-검정을 사용함.

귀무가설( $H_0$ ): 2024년과 2025년 부정 리뷰의 평균 평점에는 차이가 없다.

대립가설( $H_1$ ): 2024년과 2025년 부정 리뷰의 평균 평점에는 차이가 있다.

→ 검정통계량: -0.579, p값: 0.5638

⇒ p값이 0.05 이상이므로 귀무가설을 채택한다.

⇒ 두 연도 간 부정 리뷰의 평가에는 유의미한 차이가 없다고 본다.

🔴 [검정 c] 긍정 리뷰 표현에 따른 평점 차이

[제반조건]

- 두 집단은 서로 독립적이다.
  - 표본 수는 각각 15개이며, 등분산을 가정할 수 없다.
- 따라서 Welch의 t-검정을 수행함.

귀무가설( $H_0$ ): '좋다'와 '그저 그렇다'를 포함한 긍정 리뷰의 평균 평점에는 차이가 없다.

대립가설( $H_1$ ): 두 표현을 포함한 긍정 리뷰의 평균 평점에는 차이가 있다.

→ 검정통계량: -0.764, p값: 0.4515

⇒ p값이 0.05 이상이므로 귀무가설을 채택한다.

⇒ 표현 방식에 따른 긍정 평가에는 차이가 없다고 본다.