

추천 시스템의 목표: 자신도 좋아하리라 몰랐던 추천을 시스템이 발견. → 그에 맞는 콘텐츠를 가늠

추천 시스템의 유형

- ① 콘텐츠 기반 필터링
- ② 협업 필터링
  - ②-① 추천 이웃
  - ②-② 잠재 요인

### 콘텐츠 기반 필터링 추천 시스템

: 특정한 아이템을 산출할 경우, 그 아이템과 비슷한 콘텐츠를 가진...

ex) 영화 추천.

### 추천 이웃 협업 필터링

: 사용자 행동 양식만을 기반으로 추천을 수행하는 것. (협업 필터링 방식)

→ 사용자-아이템 평점 매트릭스 기반. 축적된 사용자 행동 데이터 기반으로 아직 평가하지 않은 아이템을 예측 평가.

	ITEM1	ITEM2	ITEM3
USER1			
USER2			
USER3			

최소행렬 ↑

추천 이웃 협업 필터링

- 사용자 기반
- 아이템 기반 (메트릭 협업 필터링)

같이 이용, 관심 유사도

### 잠재요인 협업 필터링

: 사용자-아이템 평점 행렬 데이터를 이용해 '잠재 요인'을 끄집어 낸다.

평점: 사용자의 장르 선호도 벡터와 영화의 장르별 특징 벡터의 곱.

### 행렬분해

평점행렬  $R (M \times N)$  → 사용자  $P (M \times K)$  + 아이템 행렬  $Q.T (K \times N)$

사용자    아이템

$$R = P * Q.T$$

→ P를 P나 Q로 어떻게 분해?

행렬분해는 주로 SVD 방식을 이용 (bnc) SVD는 m이 값이 많은 행렬에만 가능.

R 행렬은 최소행렬 ⇒ SGD (경사하강법) / ALS 이용.

## <실습>

### ① 콘텐츠 기반 필터링 실습 (TMDB 5000 영화 데이터)

- 장르 속성 기반 (영화/감독/배우) 간의 유사성 파악.

#### 1) 장르 속성 기반

· 장르 칼럼 값의 유사도를 비교 → 높은 평점을 받은 영화를 추천

#### 1. 필터링에 사용할 관측한 쿼리

→ genre, keywords 컬럼에는 리스트 내부에 여러 개의 딕셔너리로 구성. (가용필요)

· 각 row의 `literal_eval()` 함수를 이용. → `lst [dict 1, dict 2]` 로 생성

· apply lambda 식을 이용하여 필요한 값만 추출.  
'name'

#### 2. 장르 콘텐츠 유사도 측정

가장 간단한 방식) genres → 문자열로 변경 → count Vector로 토큰 벡터화 → 코사인 유사도

#### 3. 장르 유사도에 따라 영화 추천 함수 생성

→ 추천 기준이 되는 영화 제목 (ex) God Father), 추천할 개수를 입력

→ 선택된 값의 index 개수를 ndarray로 변환

→ genre, sorted\_ind 유사도 순으로 가져옴.

→ 출력 결과, 개수가 필요로 보임.

#### 4. 더 많은 추천 + 영화의 평점으로 필터링

\* but. 영화 평점은 저장된 데이터의 특이점 ↑

→ 평가 확률에 대한 가중치가 부여된 평점 방식을 이용.

$$\text{가중평점} = (V/(V+M)) * R + (M/(V+M)) * C$$

→ 이를 weighted\_vote - average로 서로 생성.

→ 장르 유사성이 높은 영화를 2배수 추천 → weighted\_vote 높은 순으로 1배수 최종 추천.

### ② 아이템 기반 최근접 이웃 협업 필터링 실습 (아이템 기반)

1. ratings 데이터를 row로는 USER, item을 column으로 본다. (pivot table 이용)

2. 가독성을 위해 movieID가 아닌 영화명 (movies 데이터) title로 변경.

→ ratings, movies를 merge → movieID가 아닌 title → pivot → NaN은 0으로 변경

3. 영화 간 유사도 산출. (코사인 유사도)

ratings matrix가 UserID 기준 행렬 데이터 → transpose를 해와 '영화' 기준 유사도.

### ③ 다항식 기반 최근접 이웃 행렬 필터링 (개연화)

개인화된 평가 추천: 아직 개인이 본 적이 없는 영화를 추천.

$$/ \text{사용자별 예측 평점}(R_{ui}) = \left\{ (\text{사용자 } u \text{의 모든 영화에 대한 실제 평점}) \cdot (\text{영화 } i \text{의 다른 모든 영화와의 코사인 유사도}) \right\} / \underbrace{\sum_j |S_{i,j}|}_{\text{장르인.}}$$

2. 예측 결과가 원래의 실제 평점과 얼마나 차이가 있는지를 확인. (MSE)

→ 사용자가 평점을 부여한 영화에 대해서만...

→ MSE를 감소시키는 방향...

3. 특정 영화의 가장 비슷한 유사도를 가진 영화에 대해서만 유사도 벡터를 적용.

→ 수렴시간 ↑, MSE 향상됨.

### ④ 행렬 분해를 이용한 잠재 요인 행렬 필터링

RMSE 값이 낮아 클수록 → SGD 기반...

4줄의 `get_rmse()`를 그대로 이용.

1. 영화 평점 행렬을 `USER-ITEM` 행렬로 생성.

2. `matrix-factorization()`을 이용해 행렬 분해.

3. `get_unseen_movies()` - 다시 이용하여 추천.

### ⑤ 파이썬 추천 시스템 패키지 - Surprise

Surprise는 추천적으로 2차 레벨의 데이터를 3차 레벨의 데이터로 변경. (원본 적용)

- 사용자 아이디 / 아이템 아이디 / 평점 순으로 된 데이터만 가능. (이 3개의 column만을 이용)

• OS 파일은 코드: 데이터를 삭제해야 한다.

• 무비컨트 데이터 형식이 아닌 OS의 경우, Reader 클래스를 미리 설정.

• pandas DataFrame에서 코드: 출력 순서 명시

### ⑥ 베이스라인 평점 (전체 평균 평점 + 사용자 평균 점수 + 아이템 평균 점수)

: 개인의 성향을 반영해 아이템 평가에 편향성 (bias) 요소를 반영

후한 사람, 전 사람.

### ⑦ 교차검증과 하이퍼 파라미터 튜닝

: `Cross-validate()` 이용.