

# 협업 필터링 영화 추천

## 학습 내용

- MovieLens 데이터 셋을 활용한다.
- 아이템 기반 영화 필터링으로 영화 추천을 수행해 본다.

## 목차

[01 데이터 불러오기](#)

[02 데이터 준비 및 특이값분해\(SVD\), 상관계수 구하기](#)

[03 유사영화를 찾아보기](#)

## 01 데이터 불러오기

[목차로 이동하기](#)

### 데이터 셋

- MovieLens 100K Dataset
- url : <https://grouplens.org/datasets/movielens/100k/> (<https://grouplens.org/datasets/movielens/100k/>)
- u.data
  - user\_id : 유저 정보
  - item\_id : 영화 정보
  - ratings : 평점 정보
  - timestamp : 시간 정보

In [1]:



```
import pandas as pd
import numpy as np
import sklearn
from sklearn.decomposition import TruncatedSVD
from IPython.display import display, Image
```

## 데이터 불러오기

In [2]:



```
display(Image(filename='data_u_data.png'))
```

u.data			
196	242	3	881250949↓
186	302	3	891717742↓
22	377	1	878887116↓
244	51	2	880606923↓
166	346	1	886397596↓
298	474	4	884182806↓
115	265	2	881171488↓
253	465	5	891628467↓
305	451	3	886324817↓
6	86	3	883603013↓
62	257	2	879372434↓
286	1014	5	879781125↓

In [3]:



```
col_name = ['user_id', 'item_id', 'rating', 'timestamp']
df = pd.read_csv('../data/ml-100k/u.data', sep='\\t', names=col_name)
print(df.shape)
df.head()
```

(100000, 4)

Out[3]:

	user_id	item_id	rating	timestamp
0	196	242	3	881250949
1	186	302	3	891717742
2	22	377	1	878887116
3	244	51	2	880606923
4	166	346	1	886397596

## 데이터 불러오기

### 데이터 정보

- 파일명 : u.u\_item
  - item\_id : 영화 정보
  - movie title : 영화 제목
  - release date : 출시일
  - video release date : 비디오 출시일
  - IMDb URL : IMDb URL 정보
  - unkonwn, ... : 기타 장르 정보

In [4]:



```
display(Image(filename='data_u_item.png'))
```



In [5]:



```
# 장르 분야
col_name = ['item_id', 'movie title', 'release date', 'video release date', 'IMDb URL',
            'unknown', 'Action', 'Adventure', 'Animation', 'Childrens',
            'Comedy', 'Crime', 'Documentary', 'Drama', 'Fantasy',
            'Film-Noir', 'Horror', 'Musical', 'Mystery', 'Romance',
            'Sci-Fi', 'Thriller', 'War', 'Western']

movies = pd.read_csv('../data/ml-100k/u.item', sep='|',
                    names=col_name, encoding='latin-1')

movie_names = movies[['item_id', 'movie title']]
c_movies_data = pd.merge(df, movie_names, on='item_id')
print(c_movies_data.shape)
c_movies_data.head()
```

(100000, 5)

Out[5]:

	user_id	item_id	rating	timestamp	movie title
0	196	242	3	881250949	Kolya (1996)
1	63	242	3	875747190	Kolya (1996)
2	226	242	5	883888671	Kolya (1996)
3	154	242	3	879138235	Kolya (1996)
4	306	242	5	876503793	Kolya (1996)

## 02 데이터 준비 및 특이값분해(SVD), 상관계수 구하기

[목차로 이동하기](#)

### 사용자-아이템 표 만들기

- 우리의 데이터에 대해 pivot를 사용해본다.
- 빈 값은 0으로 채우기
- 각 user에 대한 평점을 확인이 가능하다.

In [9]:



```
rating_crosstab = c_movies_data.pivot_table(values='rating',
                                             index='user_id',
                                             columns='movie title', fill_value=0)
print(rating_crosstab.shape)
rating_crosstab.head()
```

(943, 1664)

Out[9]:

movie title	'Til There Was You (1997)	1-900 (1994)	Dalmatians (1996)	101 Dalmatians (1996)	12 Angry Men (1957)	187 (1997)	2 Days in the Valley (1996)	20,000 Leagues Under the Sea (1954)	2001: A Space Odyssey (1968)	3 Ninjas: High Noon At Mega Mountain (1998)	39 Steps. The (1935)
user_id	<hr/>										
1	0	0		2	5	0	0	3	4	0	C
2	0	0		0	0	0	0	0	0	1	C
3	0	0		0	0	2	0	0	0	0	C
4	0	0		0	0	0	0	0	0	0	C
5	0	0		2	0	0	0	0	4	0	C

5 rows × 1664 columns

## 아이템-사용자 형태를 위해 행열 바꾸기

- 1664개의 영화
- 943명의 사용자

In [12]:



```
X = rating_crosstab.T
print(X.shape)
```

(1664, 943)

## SVD(특잇값 분해)

- 사이킷런을 활용하여 차원 축소 SVD를 수행.
- truncated SVD를 사용하여 차원 축소

In [14]:



```
SVD = TruncatedSVD(n_components=12, random_state=5)
resultant_matrix = SVD.fit_transform(X)
resultant_matrix.shape
```

Out[14]:

(1664, 12)

- 1664개 행(영화)과 잠재변수 12개의 열을 갖는 행렬 생성

## Correlation Pearson

- 피어슨 상관계수, 코사인 유사성과 같은 다양한 유사성 측정 지표를 사용할 수 있다.
- 피어슨 상관계수를 이용하여 상관 행렬을 만들어봄.
- 상관행렬은 단순히 상관관계를 표시하는 행렬이다.

In [15]:



```
### correlation matrix
corr_mat = np.corrcoef(resultant_matrix)
print( corr_mat.shape )
corr_mat
```

(1664, 1664)

Out[15]:

```
array([[ 1.          , -0.11573577,  0.51362284, ...,  0.38310045,
        0.20193733,  0.5065142 ],
       [-0.11573577,  1.          ,  0.05820808, ...,  0.15805829,
        0.51795357,  0.27104818],
       [ 0.51362284,  0.05820808,  1.          , ...,  0.76575655,
        0.43824619,  0.19507139],
       ...,
       [ 0.38310045,  0.15805829,  0.76575655, ...,  1.          ,
        0.18043708,  0.12115972],
       [ 0.20193733,  0.51795357,  0.43824619, ...,  0.18043708,
        1.          ,  0.20126072],
       [ 0.5065142 ,  0.27104818,  0.19507139, ...,  0.12115972,
        0.20126072,  1.          ]])
```

## 03 유사영화를 찾아보기

[목차로 이동하기](#)

### Similar Movies to Star Wars (1977)

In [16]:



```
rating_crosstab.columns.get_loc("Star Wars (1977)")
```

Out[16]:

1398

In [17]:



```
col_idx = rating_crosstab.columns.get_loc("Star Wars (1977)")
corr_specific = corr_mat[col_idx] # Star Wars (1977)의 위치 행 획득
print(corr_specific.shape)
```

(1664,)

In [18]:



```
result = pd.DataFrame({'corr_specific':corr_specific, 'Movies': rating_crosstab.columns})
print(result.shape)
result.head()
```

(1664, 2)

Out[18]:

	<b>corr_specific</b>	<b>Movies</b>
0	0.357238	'Til There Was You (1997)
1	0.421507	1-900 (1994)
2	0.593815	101 Dalmatians (1996)
3	0.722361	12 Angry Men (1957)
4	0.325221	187 (1997)

## 10개의 영화 추천

In [19]:



```
result.sort_values('corr_specific', ascending=False).head(10)
```

Out[19]:

	<b>corr_specific</b>	<b>Movies</b>
<b>1398</b>	1.000000	Star Wars (1977)
<b>1234</b>	0.988052	Return of the Jedi (1983)
<b>1460</b>	0.942655	Terminator 2: Judgment Day (1991)
<b>1523</b>	0.933978	Toy Story (1995)
<b>1461</b>	0.931701	Terminator, The (1984)
<b>1205</b>	0.925185	Raiders of the Lost Ark (1981)
<b>456</b>	0.923562	Empire Strikes Back, The (1980)
<b>570</b>	0.915965	Fugitive, The (1993)
<b>414</b>	0.914299	Die Hard (1988)
<b>44</b>	0.892894	Aliens (1986)

## (실습) Godfather, The (1972)에 대한 10개의 영화 추천해 보기

In [20]:



```
col_idx = rating_crosstab.columns.get_loc("Godfather, The (1972)")
corr_specific = corr_mat[col_idx] # Godfather, The (1972)의 위치 행 획득
print(corr_specific.shape)
```

(1664,)

In [21]:



```
result = pd.DataFrame({'corr_specific':corr_specific, 'Movies': rating_crosstab.columns})
result.sort_values('corr_specific', ascending=False).head(10)
```

Out[21]:

	<b>corr_specific</b>	<b>Movies</b>
<b>612</b>	1.000000	Godfather, The (1972)
<b>613</b>	0.921444	Godfather: Part II, The (1974)
<b>498</b>	0.921420	Fargo (1996)
<b>623</b>	0.900758	GoodFellas (1990)
<b>237</b>	0.865385	Bronx Tale, A (1993)
<b>1398</b>	0.865148	Star Wars (1977)
<b>209</b>	0.864269	Boot, Das (1981)
<b>389</b>	0.857308	Dead Man Walking (1995)
<b>622</b>	0.845558	Good, The Bad and The Ugly, The (1966)
<b>1190</b>	0.842705	Pulp Fiction (1994)

- 우리는 Godfather의 영화를 좋아하는 사람이 있다면 Godfather: Part II, Star Wars (1977)를 볼 것을 제안할 수 있다.
- 역으로 Godfather의 영화를 피하는 사람이라면 Godfather: Part II, Star Wars (1977) 를 피할 것을 제안할 수 있다.

## 실습해 보기

- Pulp Fiction (1994) 에 대한 유사 영화 10개를 추천해 보자.

## History

- 2022-11 ver 02