

[Get started](#)[Open in app](#)[Follow](#)

569K Followers



You have **2** free member-only stories left this month.

[Sign up for Medium and get an extra one](#)

Two examples of a content-based recommendation system

Content-based, weighted content-based, Numpy functions



Xue Wang Mar 26 · 4 min read ★

[Get started](#)[Open in app](#)

Photo from [Michal Matlon on Unsplash](#)

Today I would like to discuss two examples for content-based recommendation systems and some efficient array functions I learn from them. The two examples are

1: Based on item content recommendation

2: Based on weighted content recommendation

I use a simple movie set as an example and would like to focus on the main process and ignore other processes and special cases. Let's get started.

Datasets preparing:

Use the below codes to generate two datasets: `movie_df` and `review_df`

```
1 movie_id=[1,2,3,4,5]
2 user_id=[100,100,200,300, 400]
3 movie_title=['toy story', 'superman', 'titannic', 'follow me', 'minari']
```

[Get started](#)[Open in app](#)

```
6 movies={'movie_id':movie_id, 'movie_title':movie_title, 'genres':genres}
7 reviews={'movie_id':movie_id, 'user_id':user_id, 'rating':rating}
8 movies_df= pd.DataFrame(movies,columns=['movie_id', 'movie_title','gen
9 reviews_df= pd.DataFrame(reviews,columns=['movie_id', 'user_id', 'ratin
```

THE TWO TABLES AS:

movies_df

	movie_id	movie_title	genres
0	1	toy story	children,comedy
1	2	superman	drama,children
2	3	titannic	drama,romantic
3	4	follow me	horror
4	5	minari	drama

reviews_df

	movie_id	user_id	rating
0	1	100	3
1	2	100	2
2	3	200	3
3	4	300	2
4	5	400	2

Method 1: based on the movie content, make a recommendation for each user when the similarity of the content is greater than 0.

[Get started](#)[Open in app](#)

recommendation system technique, we will use the One Hot Encoding technique to convert the list of genres to a vector where each column corresponds to one possible value of the feature.

Every genre is separated by a “,” so we simply have to call the split function on

```
1 movies_df['genres'] = movies_df.genres.str.split(',')
2
3 movies_genres_df = movies_df.copy()
4 #For every row in the dataframe, iterate through the list of genres and
5 for index, row in movies_df.iterrows():
6     for genre in row['genres']:
7         movies_genres_df.at[index, genre] = 1
8
9 movies_genres_df.head()
```

	movie_id	movie_title	genres	children	comedy	drama	romantic	horror
0	1	toy story	[children, comedy]	1.0	1.0	0.0	0.0	0.0
1	2	superman	[drama, children]	1.0	0.0	1.0	0.0	0.0
2	3	titannic	[drama, romantic]	0.0	0.0	1.0	1.0	0.0
3	4	follow me	[horror]	0.0	0.0	0.0	0.0	1.0
4	5	minari	[drama]	0.0	0.0	1.0	0.0	0.0

Step 1: Calculate the movie-movie similarity matrix :

```
1 # Subset using the dummy variables
2 movie_genres = np.array(movies_genres_df.iloc[:,3:])
3 # dot product to obtain a movie x movie matrix of similarities
4 dot_prod_movies_genres = movie_genres .dot(np.transpose(movie_genres ))
```

[Get started](#)[Open in app](#)

```
[1., 2., 1., 0., 1.],
[0., 1., 2., 0., 1.],
[0., 0., 0., 1., 0.],
[0., 1., 1., 0., 1.]])
```

The dot product shows the similarity among the movies.

Step 2: Find similar movies: here the standard is that if there is a relationship then choose to recommend as the small dataset.

```
1 def find_similar_movies(movie_id):
2     '''
3     INPUT
4     movie_id - a movie_id
5     OUTPUT
6     similar_movies - an array of the most similar movies by title
7     '''
8     # find the row of each movie id
9     movie_idx = np.where(movies_genres_df['movie_id'] == movie_id)[0][
10
11
12     # find the most similar movie indices - to start I said they need
13
14     similar_idx = np.where(dot_prod_movies_genres[movie_idx] >= 1)[0]
15
16
17     # pull the movie titles based on the indices
18     similar_movies = np.array(movies_genres_df.iloc[similar_idx, ]['m
```

which only choose the highest value of the similarity.

```
1 # only choose the most similar movie
```



Get started

Open in app

np.where() shows the item position (index) in the table.

Make a test for movie_id = 1:

```
find_similar_movies(1)
array(['toy story', 'superman'], dtype=object)
```

A similar movie with toy story is superman, whose genre is children.
The result is as expected.

Then get a list of movie name, if the name of the movie is a list:

```
1 def get_movie_names(movie_ids):
2     '''
3     INPUT
4     movie_ids - a list of movie_ids
5     OUTPUT
6     movies - a list of movie names associated with the movie_ids
7     '''
8     movie_lst = list(movies_genres_df[movies_genres_df['movie_id'].isin
9     return movie_lst
```

Step 6: Make recommendations for a specific user.

```
1 def make_recs(user_id):
2
3     recs=np.array([])
4     # Pull only the reviews the user has seen
5     reviews_temp = reviews_df[reviews_df['user_id'] == user_id]
6     movies_temp = np.array(reviews_temp['movie_id'])
```

[Get started](#)[Open in app](#)

```
10     temp_recs = np.setdiff1d(rec_movies, movie_names)
11     recs=list(temp_recs)
12
```

```
make_recs(100)
['minari', 'titannic']
```

Both of them are drama, as user 100 has also reviewed drama ‘superman’, so the recommendation makes sense.

Method 2: Based on weighted content

Method 1 is easy to understand, but it seems the rating information has not been used. Now I want to integrate the info to calculate the weighted genres.

For example, I want to construct weighted genres based on the user’s rating. Let’s choose user_id=100.

Step 1: Filter out the movies which are rated by user 100 and get the genres only:

```
1  #Filtering out the movies from the review
2  movie_id_df=reviews_df.loc[reviews_df['user_id']==user_id]
3  user_movies = movies_genres_df[movies_genres_df['movie_id'].isin(movie_
4  user_movies = user_movies.iloc[:,3:]
```

11.00.01.00.0

Step 2: Get weighted genres for this user:

```
1 rating_df=reviews_df.loc[reviews_df['user_id']==user_id]['rating']
2 #Dot produt to get weights
3 userProfile = userMovies.transpose().dot(rating_df)
```

get weighted movie genres by heated with ❤️ by CitiHubview raw

```
children    5.0
comedy      3.0
drama       2.0
romantic    0.0
horror      0.0
dtype: float64
```

Step 3: Get a recommendation (user-reviewed also included)

```
1 recommendation_array=movie_genres.transpose().dot(userProfile)/(userPro
2 recommendation_series=pd.Series(recommendation_array).sort_values(ascen
3 recommendation_keys=[keys for keys, value in recommendation_series.item
4 rec_movies=movies_df.loc[movies_df['movie_id'].isin(movie_ids)]
```

view raw

	movie_id	movie_title	genres
0	1	toy story	[children, comedy]
1	2	superman	[drama, children]
2	3	titannic	[drama, romantic]

[Get started](#)[Open in app](#)

the user.

Step 4: Filter the reviewed movies

```
1 user_reviews = reviews_df[reviews_df['user_id'] == 100]
2 movies_reviews = np.array(user_reviews['movie_id'])
3 movie_names = np.array(get_movie_names(movies_reviews))
4 rec_movies=np.array(rec_movies['movie_title'])
5 recs=np.setdiff1d(rec_movies, movie_names)
```

```
recs
```

```
array(['titannic'], dtype=object)
```

In this case, as minari's weighted genres for user 100 is 0.2, so it is excluded from the recommendation.

Summary:

- From the two cases, we can see that no reviewed genres will not be recommended, which is the character of a content-based recommendation system. It is highly personalized for the user.
- To some degree, a recommendation system is like an art, and you can also create your criteria to adjust to your target, like in this story I have displayed various criteria even for the same method.
- There are some array functions which I believe are very efficient when dealing with recommendation system, and would like to summarize here again:

[Get started](#)[Open in app](#)

and use it.

- 2) `np.dot(np.transpose())`: this dot product can be used not only for the item itself but also for getting the weighted items.
- 3) `np.where()`: here I use it to find the location(index) of the item.
- Of course, you can also extrapolate them to other situations, not only the recommendation system. Because of the limitation, I haven't enough time to show how the functions are efficient, and I explain them [in this story](#) in detail.

Thank you for your reading.

Sign up for The Variable

By Towards Data Science

Every Thursday, the Variable delivers the very best of Towards Data Science: from hands-on tutorials and cutting-edge research to original features you don't want to miss. [Take a look.](#)

[Get this newsletter](#)[Recommendation System](#)[Arrays](#)[Python](#)[Content Based Filtering](#)



Get started

Open in app

Get the Medium app

