

TRƯỜNG ĐẠI HỌC SƯ PHẠM KỸ THUẬT TP. HỒ CHÍ MINH
KHOA ĐÀO TẠO CHẤT LƯỢNG CAO
NGÀNH CÔNG NGHỆ THÔNG TIN

☺ ☹ ☹



BÁO CÁO TIẾN ĐỘ -TUẦN 05
MÔN: ĐỒ ÁN 3
ĐỀ TÀI: TÌM HIỂU VỀ THUẬT TOÁN
RECOMMENDATION

GVHD : Thầy Huỳnh Xuân Phụng

SVTH :

Nguyễn Thành Như

17110202

Võ Ngọc Thuận

17110234

TP. Hồ Chí Minh, tháng 11 năm 2020

Trần

1. Demo Recommendation System with Content-based Filtering

1.1. Code Demo

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.feature_extraction.text import CountVectorizer
4 from sklearn.metrics.pairwise import cosine_similarity
5 ##### helper functions. Use them when needed #####
6 def get_title_from_index(index):
7     return df[df.index == index]["title"].values[0]
8
9 def get_index_from_title(title):
10    return df[df.title == title]["index"].values[0]
11 #####
12
13 ##Step 1: Read CSV File
14 df = pd.read_csv("movie_dataset.csv")
15 print (df.columns)
16
17 ##Step 2: Select Features
18 features = ['keywords','cast','genres','director']
19
20 ##Step 3: Create a column in DF which combines all selected features
21 for feature in features:
22     df[feature] = df[feature].fillna('')
23
24 def combine_features(row):
25     try:
26         return row['keywords']+" "+row['cast']+" "+row['genres']+" "+row['director']
27     except:
28         print ("Error:", row)
29
30 df["combined_features"] = df.apply(combine_features,axis=1)
31 df.iloc[0].combined_features
32 ##Step 4: Create count matrix from this new combined column
33 cv = CountVectorizer()
34 count_matrix = cv.fit_transform(df["combined_features"])
35
36 ##Step 5: Compute the Cosine Similarity based on the count_matrix
37 cosine_sim = cosine_similarity(count_matrix)
38
```

```

39 ## Step 6: Get index of this movie from its title
40 def get_title_from_index(index):
41     return df[df.index == index]["title"].values[0]
42 def get_index_from_title(title):
43     return df[df.title == title]["index"].values[0]
44 movie_user_likes = "Avatar"
45 movie_index = get_index_from_title(movie_user_likes)
46 similar_movies = list(enumerate(cosine_sim[movie_index]))
47
48 ## Step 7: Get a list of similar movies in descending order of similarity score
49 sorted_similar_movies = sorted(similar_movies, key=lambda x:x[1], reverse=True)[1:]
50
51 ## Step 8: Print titles of first 50 movies
52 i=0
53 print("Top 5 similar movies to "+movie_user_likes+" are:\n")
54 for element in sorted_similar_movies:
55     print(get_title_from_index(element[0]))
56     i=i+1
57     if i>5:
58         break
59

```

1.2. Giải thích code:

- **import pandas as pd** : chúng ta tiến hành import thư viện Pandas, nó là một thư viện Python cung cấp các cấu trúc dữ liệu nhanh, mạnh mẽ, linh hoạt và mang hàm ý. Chúng em dùng nó để đọc dữ liệu từ file csv.
- **from scipy import sparse** : scipy là một thư viện mã nguồn mở các thuật toán và các công cụ toán học cho Python, được xây dựng trên các đối tượng mảng NumPy tạo thành ngăn xếp NumPy bao gồm các công cụ như Pandas, SymPy và Matplotlib. SciPy cung cấp khá nhiều module tính toán từ đại số tuyến tính, tích phân, vi phân, nội suy đến xử lý ảnh, fourier transform... và cụ thể ở đây ta import gói ma trận sparse và các đoạn chương trình liên quan
- **from sklearn.metrics.pairwise import cosine_similarity** : cosine_similarity được lấy từ sklearn.metrics.pairwise để tính toán độ tương tự giữa 2 vectors
- **Step 1: df=pd.read_csv("movie_dataset.csv")**: gán biến df và đọc biến df với dữ liệu là file movie_dataset.csv bằng việc đọc vào một file .csv bằng cách sử dụng hàm read_csv và được trả về 1 dataframe.

```

13 ##Step 1: Read CSV File
14 df = pd.read_csv("movie_dataset.csv")
15 print(df.columns)
16

```

- **Step 2: features = ['keywords','cast','genres','director']**: chọn 4 features 'keywords', 'cast', 'genres', 'director' làm bộ tính năng cho nội dung của movie_dataset

```

17 ##Step 2: Select Features
18 features = ['keywords', 'cast', 'genres', 'director']
19

```

- **Step 3:** Nhiệm vụ tiếp theo của chúng ta là tạo một hàm để kết hợp các giá trị của các cột này thành một chuỗi duy nhất. chúng ta cần gọi hàm này trên mỗi hàng của khung dữ liệu. Tuy nhiên, trước khi làm điều đó, chúng ta cần clean và xử lý trước dữ liệu để sử dụng. Điền tất cả các giá trị NaN bằng chuỗi trống trong khung dữ liệu.

```
20 ##Step 3: Create a column in DF which combines all selected features
21 for feature in features:
22     df[feature] = df[feature].fillna('')
23
24 def combine_features(row):
25     try:
26         return row['keywords']+" "+row['cast']+" "+row['genres']+" "+row['director']
27     except:
28         print ("Error:", row)
29
30 df["combined_features"] = df.apply(combine_features,axis=1)
31 df.iloc[0].combined_features
```

- **Step 4:** `cv = CountVectorizer()`: tạo object CountVectorizer().
`count_matrix = cv.fit_transform(df["combined_features"]):`
 combine chuỗi string(movie contents) với object CountVectorizer().

```
32 ##Step 4: Create count matrix from this new combined column
33 cv = CountVectorizer()
34 count_matrix = cv.fit_transform(df["combined_features"])
```

- **Step 5:** Chúng ta cần lấy ma trận tương tự cosine từ ma trận đếm.

```
36 ##Step 5: Compute the Cosine Similarity based on the count_matrix
37 cosine_sim = cosine_similarity(count_matrix)
```

- **Step 6:** Xác định hai hàm trợ giúp để lấy title phim từ chỉ mục phim và ngược lại.

```
39 ## Step 6: Get index of this movie from its title
40 def get_title_from_index(index):
41     return df[df.index == index]["title"].values[0]
42 def get_index_from_title(title):
43     return df[df.title == title]["index"].values[0]
44 movie_user_likes = "Avatar"
45 movie_index = get_index_from_title(movie_user_likes)
46 similar_movies = list(enumerate(cosine_sim[movie_index]))
```

Bước tiếp theo là lấy tên phim mà người dùng đang thích. Sau đó, tìm chỉ số của bộ phim đó. Sau đó, truy cập vào hàng tương ứng với bộ phim này trong ma trận tương tự. Như vậy, chúng ta sẽ nhận được điểm tương đồng của tất cả

các phim khác từ phim hiện tại. Sau đó, thống kê thông qua tất cả các điểm tương tự của bộ phim đó để tạo ra một bộ chỉ số phim và điểm tương tự. Điều này sẽ chuyển đổi một hàng điểm số tương tự như thế này- [1 0,5 0,2 0,9] thành hàng này- [(0, 1) (1, 0,5) (2, 0,2) (3, 0,9)].

```
45 movie_user_likes = "Avatar"
46 movie_index = get_index_from_title(movie_user_likes)
47 similar_movies = list(enumerate(cosine_sim[movie_index]))
```

- **Step 7:** Sắp xếp danh sách similar_movies theo điểm tương tự theo thứ tự giảm dần. Vì phim tương tự nhất với một phim nhất định sẽ là chính nó, loại bỏ phần tử đầu tiên sau khi sắp xếp các phim.

```
49 ## Step 7: Get a list of similar movies in descending order of similarity score
50 sorted_similar_movies = sorted(similar_movies, key=lambda x:x[1], reverse=True)[1:]
```

- **Step 8:** Chạy vòng lặp để in 5 entries đầu tiên từ danh sách sorted_similar_movies.

```
49 ## Step 7: Get a list of similar movies in descending order of similarity score
50 sorted_similar_movies = sorted(similar_movies, key=lambda x:x[1], reverse=True)[1:]
51
52 ## Step 8: Print titles of first 50 movies
53 i=0
54 print("Top 5 similar movies to "+movie_user_likes+" are:\n")
55 for element in sorted_similar_movies:
56     print(get_title_from_index(element[0]))
57     i=i+1
58     if i>5:
59         break
```

Kết quả: Top 5 bộ phim liên quan, gần với phim Avatar nhất là:

```
23
24 def combine_features(row):
25     try:
26         return row['keywords']+" "+row['cast']+" "+row['genres']+" "+row['director']
27     except:
28         print ("Error:", row)
29
30 df["combined_features"] = df.apply(combine_features,axis=1)
31 df.iloc[0].combined_features
32
33 ##Step 4: Create count matrix from this new combined column
34 cv = CountVectorizer()
35 count_matrix = cv.fit_transform(df["combined_features"])
36
37 ##Step 5: Compute the Cosine Similarity based on the count_matrix
38 cosine_sim = cosine_similarity(count_matrix)
39
40 ## Step 6: Get index of this movie from its title
41 def get_title_from_index(index):
42     return df[df.index == index]["title"].values[0]
43
44 def get_index_from_title(title):
45     return df[df.title == title]["index"].values[0]
46
47 movie_user_likes = "Avatar"
48 movie_index = get_index_from_title(movie_user_likes)
49 similar_movies = list(enumerate(cosine_sim[movie_index]))
50 sorted_similar_movies = sorted(similar_movies,key=lambda x:x[1],reverse=True)[1:]
51
52 ## Step 7: Get a list of similar movies in descending order of similarity score
53 i=0
54 print("Top 5 similar movies to "+movie_user_likes+" are:\n")
55 for element in sorted_similar_movies:
56     print(get_title_from_index(element[0]))
57     i=i+1
58     if i>5:
59         break
60
```

Variable explorer

Name	Type	Size	Value
cosine_sim	float64	(4803, 4803)	[[1. 0.10540926 0.12038585 ... 0. 0. 0. ...
df	DataFrame	(4803, 25)	Column names: index, budget, genres, homepage, id, keywords, original_ ...
element	tuple	2	(3158, 0.3333333333333333)
feature	str	1	director
features	list	4	['keywords', 'cast', 'genres', 'director']
i	int	1	6

IPython console

```
IPython 7.4.0 -- An enhanced Interactive Python.
In [1]: runfile('E:/DoAn3/movie_recommender/movie_recommender.py', wdir='E:/DoAn3/
movie_recommender')
Index(['index', 'budget', 'genres', 'homepage', 'id', 'keywords',
'original_language', 'original_title', 'overview', 'popularity',
'production_companies', 'production_countries', 'release_date',
'runtime', 'spoken_languages', 'status', 'tagline', 'title',
'vote_average', 'vote_count', 'cast', 'crew', 'director'],
dtype='object')
Top 5 similar movies to Avatar are:
Guardians of the Galaxy
Aliens
Star Wars: Clone Wars: Volume 1
Star Trek Into Darkness
Star Trek Beyond
Alien
```

Top 5 bộ phim liên quan, gần với bộ phim Avatar nhất là:

Guardians of the Galaxy
Aliens
Star Wars: Clone Wars: Volume 1
Star Trek Into Darkness
Star Trek Beyond
Alien

II. Công việc tuần sau:

- Demo dựa theo Clustering User.