

# HỆ KHUYẾN NGHỊ

## BÀI TẬP THỰC HÀNH TUẦN 3

### KHUYẾN NGHỊ DỰA TRÊN NỘI DUNG

#### 1. Quy định về việc nộp bài

- Thời gian: Được giảng viên thiết lập trên hệ thống Moodle.
- Hình thức nộp: Trên Moodle.
- Bài nộp bao gồm các file **.ipynb** trong một folder và nén lại thành một tập tin (**.zip**).
- Cách đặt tên: **BTTH3\_MSSV.zip**
- Công cụ thực hành. **Google colab**
- Lưu ý: Sai quy định thì sẽ nhận 0 điểm.

#### 2. Nội dung thực hành

##### 2.1. Theo dõi giảng viên hướng dẫn

DL mẫu:

```
train_set = ["Bộ phim kể về cuộc chiến <html> đấu chống lại <html> những  
linh hồn ác quỷ.",  
            "Cuộc đột nhập # $ của Thom một kỹ sư thiên tài vào một ngân  
hàng bí ẩn tại Hồng Kông.",  
            "Phim về những rắc rối hài hước đời thường và thế giới tình yêu  
trong độ tuổi 30.",  
            "Về các siêu <url> anh hùng cùng hợp tác và (chống) # lại mỗi  
nguy hiểm mang quy mô quốc tế."]  
  
test_set = ["Phim kể về tình yêu hài hước và quá trình sống chung rắc rối  
giữa một biên kịch với một diễn viên nam."]
```

- Tiền xử lý dữ liệu
  - Tách từ trên tiếng Anh và tiếng Việt
  - Loại bỏ stopwords từ trên tiếng Anh và tiếng Việt  
Link stopwords tiếng Việt: <https://www.kaggle.com/mpwolke/vietnamese-stopwords-w2v/notebook>
  - Chuẩn hóa từ

- Chuyển chuỗi về hoa, thường, xóa các thẻ HTML, dấu câu, ... từ biểu thức chính quy (Regular Expression)
- Tùy vào miền dữ liệu chúng ta có thể xử lý emoji.
- Xây dựng ma trận vector TF-IDF
- Khuyến nghị sản phẩm cho người dùng nhờ vào độ đo Cosine

## 2.2. Yêu cầu về nhà

Dựa trên overview của các bộ phim trong bộ dữ liệu “The Movies Dataset” cài đặt thuật toán khuyến nghị dựa trên dung.

- Link download: <https://www.kaggle.com/rounakbanik/the-movies-dataset>

	A	B	C	D	E	F	G	H	I	J	K	
	adult	belongs_to_collection	budget	genres	homepage	id	imdb_id	original_language	original_title	overview	popularity	poster_path
1	FALSE	["id": 10194, "name": "Toy Story"]	30000000	["id": 16, "name": "Adventure"]	http://toystory.com	862	tt0114709	en	Toy Story	Led by Woody, Andy's toys live happily in his room until Andy's birthday brings new toys into town.	21.946.943	/rhlRbceOE9IR
2	FALSE	["id": 119050, "name": "Grumpy Old Men"]	65000000	["id": 12, "name": "Comedy"]	http://www.grumpyoldmen.com	8844	tt0113497	en	Grumpy Old Men	When siblings Judy and Peter discover an enchanted board game that opens the door to a magical world, it's up to the grumpy old men to save the day.	17.015.539	/vzmL6fP7aPK
3	FALSE	["id": 10749, "name": "Roman Holiday"]	0	["id": 10749, "name": "Roman Holiday"]	http://www.roman-holiday.com	15602	tt0113228	en	Grumpy Old Men	A family wedding reignites the ancient feud between next-door neighbors and bitter adversaries.	117.129	/6ksm1sjKMfI
4	FALSE	["id": 96871, "name": "Father of the Bride Part II"]	16000000	["id": 35, "name": "Comedy"]	http://www.fatherofthebride.com	31357	tt0114885	en	Father of the Bride Part II	Cheated on, mistreated and stepped on, the women are holding their breath waiting to exhale.	3.859.495	/16KOMpEaLV
5	FALSE	["id": 96871, "name": "Father of the Bride Part II"]	0	["id": 35, "name": "Comedy"]	http://www.fatherofthebride.com	11862	tt0113041	en	Father of the Bride Part II	Just when George Banks has recovered from his daughter's wedding, he receives word that she is getting married again.	8.387.519	/e64sOI48hQX
6	FALSE	["id": 96871, "name": "Father of the Bride Part II"]	60000000	["id": 28, "name": "Action"]	http://www.fatherofthebride.com	949	tt0113277	en	Heat	Obsessive master thief, Neil McCauley leads a top-notch crew on various insane heists.	17.924.927	/zMyfPUelumi
7	FALSE	["id": 96871, "name": "Father of the Bride Part II"]	58000000	["id": 35, "name": "Comedy"]	http://www.fatherofthebride.com	11860	tt0114319	en	Sabrina	An ugly duckling having undergone a remarkable change, still harbors feelings for her prince.	6.677.277	/jQh15ySf87b
8	FALSE	["id": 96871, "name": "Father of the Bride Part II"]	0	["id": 28, "name": "Action"]	http://www.fatherofthebride.com	45325	tt0112302	en	Tom and Huck	A mischievous young boy, Tom Sawyer, witnesses a murder by the deadly Injun.	2.561.161	/vG05Qa55p7
9	FALSE	["id": 96871, "name": "Father of the Bride Part II"]	35000000	["id": 28, "name": "Action"]	http://www.fatherofthebride.com	9091	tt0114576	en	Sudden Death	International action superstar Jean Claude Van Damme teams with Powers of Ten.	523.158	/eWwK060T5
10	FALSE	["id": 645, "name": "James Bond: The World Is Not Enough"]	58000000	["id": 12, "name": "Action"]	http://www.bond.com	710	tt0113189	en	GoldenEye	James Bond must unmask the mysterious head of the Janus Syndicate and prevent a global nuclear war.	14.686.036	/5c0ovj741Kn
11	FALSE	["id": 645, "name": "James Bond: The World Is Not Enough"]	62000000	["id": 35, "name": "Comedy"]	http://www.bond.com	9087	tt0112346	en	The American President	Widowed U.S. president Andrew Shepherd, one of the world's most powerful men, falls for a woman who is the daughter of a man who killed his father.	6.318.445	/jymPNGLzGPI
12	FALSE	["id": 645, "name": "James Bond: The World Is Not Enough"]	0	["id": 35, "name": "Comedy"]	http://www.bond.com	12110	tt0112896	en	Dracula: Dead and Loving It	When a lawyer shows up at the vampire's doorstep, he falls prey to his charms.	5.430.331	/jve4CgYfYtN0
13	FALSE	["id": 117693, "name": "Baltico"]	0	["id": 10751, "name": "Family"]	http://www.baltico.com	21032	tt0112453	en	Balto	An outcast half-wolf risks his life to prevent a deadly epidemic from ravaging his town.	12.140.733	/gVSPCAVCVN
14	FALSE	["id": 117693, "name": "Baltico"]	44000000	["id": 36, "name": "History"]	http://www.baltico.com	10858	tt0113987	en	Nixon	An all-star cast powers this epic look at American President Richard M. Nixon.	5.092	/c1CkmXEIRh
15	FALSE	["id": 117693, "name": "Baltico"]	98000000	["id": 28, "name": "Action"]	http://www.baltico.com	1408	tt0112760	en	Cutthroat Island	Morgan Adams and her slave, William Shaw, are on a quest to recover the treasure of a lost civilization.	7.284.477	/odM9973kv9
16	FALSE	["id": 117693, "name": "Baltico"]	52000000	["id": 18, "name": "Drama"]	http://www.baltico.com	524	tt0112641	en	Casino	The life of the gambling paradise, Las Vegas, and its dark mafia underworld.	10.137.389	/xos17ibXBdF
17	FALSE	["id": 117693, "name": "Baltico"]	16500000	["id": 18, "name": "Drama"]	http://www.baltico.com	4584	tt0114388	en	Sense and Sensibility	Rich Mr. Dashwood dies, leaving his second wife and her daughters poor by the wayside.	10.673.167	/IA9HTyB48b6
18	FALSE	["id": 117693, "name": "Baltico"]	4000000	["id": 80, "name": "Crime"]	http://www.baltico.com	5	tt0113101	en	Four Rooms	It's Ted the Bellhop's first night on the job...and the hotel's very unusual guests.	9.026.586	/eQs5h8nrkxI
19	FALSE	["id": 3167, "name": "Ace Ventura: When Nature Calls"]	30000000	["id": 80, "name": "Crime"]	http://www.aceventura.com	9273	tt0112881	en	Ace Ventura: When Nature Calls	Ace Ventura: When Nature Calls	8.205.448	/wRlGnJhExcc
20	FALSE	["id": 3167, "name": "Ace Ventura: When Nature Calls"]	60000000	["id": 28, "name": "Action"]	http://www.aceventura.com	11517	tt0113845	en	Money Train	A vengeful New York transit cop decides to steal a trainload of subway fares.	7.337.906	/jSozzV0R2kd
21	FALSE	["id": 91698, "name": "Chili"]	30250000	["id": 35, "name": "Comedy"]	http://www.chili.com	8012	tt0113161	en	Get Shorty	Chili Palmer is a Miami mobster who gets sent by his boss, the psychopathic	12.669.608	/VWV0UuGQA
22	FALSE	["id": 91698, "name": "Chili"]	0	["id": 18, "name": "Drama"]	http://www.chili.com	1710	tt0112722	en	Copycat	An agoraphobic psychologist and a female detective must work together to find a serial killer.	10.701.801	/80ceG5oIk2
23	FALSE	["id": 91698, "name": "Chili"]	50000000	["id": 28, "name": "Action"]	http://www.chili.com	9691	tt0112401	en	Assassins	Assassin Robert Rath arrives at a funeral to kill a prominent mobster, only to find out he's the son of the mobster he was hired to kill.	11.065.939	/AAKMP7Dg4
24	FALSE	["id": 91698, "name": "Chili"]	0	["id": 18, "name": "Drama"]	http://www.chili.com	12665	tt0114168	en	Powder	Harassed by classmates who won't accept his shocking appearance, a shy yet brilliant boy is recruited by a	12.133.094	/JuRkxOCQg
25	FALSE	["id": 91698, "name": "Chili"]	3600000	["id": 18, "name": "Drama"]	http://www.chili.com	451	tt0113627	en	Leaving Las Vegas	Ben Sanderson, an alcoholic Hollywood screenwriter who lost everything because of his drinking, is recruited by	10.332.025	/37aHJm5h5
26	FALSE	["id": 91698, "name": "Chili"]	0	["id": 18, "name": "Drama"]	http://www.chili.com	16420	tt0114057	en	Othello	The evil Iago pretends to be friend of Othello in order to manipulate him to kill Desdemona.	1.845.899	/jQMBXECQm
27	FALSE	["id": 91698, "name": "Chili"]	12000000	["id": 35, "name": "Comedy"]	http://www.chili.com	9263	tt0114011	en	Now and Then	Waxing nostalgic about the bittersweet passage from childhood to puberty, a group of girls	8.681.325	/jQMBXECQm
28	FALSE	["id": 91698, "name": "Chili"]	0	["id": 18, "name": "Drama"]	http://www.chili.com	17015	tt0114117	en	Persuasion	This film adaptation of Jane Austen's last novel follows Anne Elliot, the daughter of a wealthy	2.828.434	/eV06eqw4ld
29	FALSE	["id": 91698, "name": "Chili"]	18000000	["id": 14, "name": "Fantasy"]	http://www.chili.com	902	tt0112682	fr	La Cité des Enfants Perdus	A scientist in a surrealist society kidnaps children to steal their dreams, hoping to create a	9.822.423	/eV06eqw4ld
30	FALSE	["id": 91698, "name": "Chili"]	0	["id": 18, "name": "Drama"]	http://www.chili.com	37557	tt0115012	zh	ÊsâÄäÊsâÖpâÊsâÄä=â	A provincial boy related to a Shanghai crime family is recruited by his uncle to	1.100.915	/qCocCo7nYC
31	FALSE	["id": 91698, "name": "Chili"]	0	["id": 18, "name": "Drama"]	http://www.chili.com	9909	tt0112792	en	Dangerous Minds	Former Marine Louanne Johnson lands a gig teaching in a pilot program for inner-city	9.481.338	/j5lee3QmY0
32	FALSE	["id": 91698, "name": "Chili"]	29500000	["id": 878, "name": "Science Fiction"]	http://www.chili.com	63	tt0114746	en	Twelve Monkeys	In the year 2035, convict James Cole reluctantly volunteers to be sent back in time to	12.297.305	/65J9wDu3Yug
33	FALSE	["id": 91698, "name": "Chili"]	0	["id": 10749, "name": "Roman Holiday"]	http://www.chili.com	78803	tt0114057	fr	Gulliver's Travels	Gulliver's Travels	0.745547	/64rN8e84e

- Sử dụng file **movies\_metadata.csv**: Chứa thông tin của hơn 45.000 bộ phim có trong bộ dữ liệu Full MovieLens. Bộ dữ liệu gồm nhiều cột khác nhau, tuy nhiên chúng ta sẽ sử dụng cột “**Overview**” nội dung chính của phim để xây dựng hệ khuyến nghị của mình.
- Thử chọn ra top 10 bộ phim cho một người dùng đã từng xem một bộ phim có tên là “**Father of the Bride Part II**”