

```
1 import pandas as pd
2
3 movies = pd.read_csv("/content/movies.csv")
4
```

```
1 movies.head()
```

	movieId	title	genres
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	2	Jumanji (1995)	Adventure Children Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama Romance
4	5	Father of the Bride Part II (1995)	Comedy

```
1 import re
2
3 def clean_title(title):
4     title = re.sub("[^a-zA-Z0-9 ]", "", title)
5     return title
```

```
1 movies["clean_title"] = movies["title"].apply(clean_title)
```

```
1 movies
```

	movieId	title	genres	clean_title
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	Toy Story 1995
1	2	Jumanji (1995)	Adventure Children Fantasy	Jumanji 1995
2	3	Grumpier Old Men (1995)	Comedy Romance	Grumpier Old Men 1995
3	4	Waiting to Exhale (1995)	Comedy Drama Romance	Waiting to Exhale 1995
4	5	Father of the Bride Part II (1995)	Comedy	Father of the Bride Part II 1995
...
62418	209157	We (2018)	Drama	We 2018
62419	209159	Window of the Soul (2001)	Documentary	Window of the Soul 2001
62420	209163	Bad Poems (2018)	Comedy Drama	Bad Poems 2018
62421	209169	A Girl Thing (2001)	(no genres listed)	A Girl Thing 2001
62422	209171	Women of Devil's Island (1962)	Action Adventure Drama	Women of Devils Island 1962

62423 rows × 4 columns

```
1 from sklearn.feature_extraction.text import TfidfVectorizer
2 vectorizer = TfidfVectorizer(ngram_range=(1,2))
3
4 tfidf = vectorizer.fit_transform(movies["clean_title"])
```

```
1 from sklearn.metrics.pairwise import cosine_similarity
2 import numpy as np
3
4 def search(title):
5     title = clean_title(title)
6     query_vec = vectorizer.transform([title])
7     similarity = cosine_similarity(query_vec, tfidf).flatten()
8     indices = np.argpartition(similarity, -5)[-5:]
9     results = movies.iloc[indices].iloc[::-1]
10
11     return results
```

```
1 import ipywidgets as widgets
2 from IPython.display import display
3
4 movie_input = widgets.Text(
5     value='Toy Story',
6     description='Movie Title:',
7     disabled=False
8 )
```

```
9 movie_list = widgets.Output()  
10  
11 def on_type(data):  
12     with movie_list:  
13         movie_list.clear_output()  
14         title = data["new"]  
15         if len(title) > 5:  
16             display(search(title))  
17  
18 movie_input.observe(on_type, names='value')  
19  
20  
21 display(movie_input, movie_list)
```

Movie Title:

```
1 movie_id = 1  
2  
3 movie = movies[movies["movieId"] == movie_id]
```

```
1 ratings = pd.read_csv("/content/ratings.csv")
```

```
1 ratings.dtypes
```

	0
userId	int64
movieId	int64
rating	float64
timestamp	int64

dtype: object

```
1 similar_users = ratings[(ratings["movieId"] == movie_id) & (ratings["rating"] > 4)]["userId"].unique()
```

```
1 similar_user_recbs = ratings[(ratings["userId"].isin(similar_users)) & (ratings["rating"] > 4)]["movieId"]
```

```
1 similar_user_recbs = similar_user_recbs.value_counts() / len(similar_users)  
2  
3 similar_user_recbs = similar_user_recbs[similar_user_recbs > .10]
```

```
1 all_users = ratings[(ratings["movieId"].isin(similar_user_recbs.index)) & (ratings["rating"] > 4)]
```

```
1 all_user_recbs = all_users["movieId"].value_counts() / len(all_users["userId"].unique())
```

```
1 rec_percentages = pd.concat([similar_user_recbs, all_user_recbs], axis=1)  
2 rec_percentages.columns = ["similar", "all"]
```

```
1 rec_percentages
```

```
similar      all
movieId
1    1.000000  0.124861
318   0.444989  0.342147
260   0.400146  0.221404
356   0.369742  0.235092
296   0.367358  0.285047
...
59315  0.103862  0.053784
778    0.103266  0.074634
953    0.103067  0.046076
551    0.101146  0.041262
1225   0.100351  0.057175
```

111 rows × 2 columns

```
1 rec_percentages["score"] = rec_percentages["similar"] / rec_percentages["all"]
```

```
1 rec_percentages = rec_percentages.sort_values("score", ascending=False)
```

```
1 rec_percentages.head(10).merge(movies, left_index=True, right_on="movieId")
```

	similar	all	score	movieId	title	genres	clean_title
0	1.000000	0.124861	8.008876	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	Toy Story 1995
3021	0.281115	0.053825	5.222752	3114	Toy Story 2 (1999)	Adventure Animation Children Comedy Fantasy	Toy Story 2 1999
2264	0.111413	0.025316	4.400826	2355	Bug's Life, A (1998)	Adventure Animation Children Comedy	Bugs Life A 1998
14813	0.153011	0.035150	4.353060	78499	Toy Story 3 (2010)	Adventure Animation Children Comedy Fantasy IMAX	Toy Story 3 2010
4780	0.235212	0.070532	3.334840	4886	Monsters, Inc. (2001)	Adventure Animation Children Comedy Fantasy	Monsters Inc 2001
6258	0.227065	0.071665	3.168428	6377	Finding Nemo (2003)	Adventure Animation Children Comedy	Finding Nemo 2003
580	0.214016	0.067554	3.168056	588	Aladdin (1992)	Adventure Animation Children Comedy Musical	Aladdin 1992
					Beauty and the		Beauty and the

```
1 def find_similar_movies(movie_id):
2     similar_users = ratings[(ratings["movieId"] == movie_id) & (ratings["rating"] > 4)]["userId"].unique()
3     similar_user_recs = ratings[(ratings["userId"].isin(similar_users)) & (ratings["rating"] > 4)]["movieId"]
4     similar_user_recs = similar_user_recs.value_counts() / len(similar_users)
5
6     similar_user_recs = similar_user_recs[similar_user_recs > .10]
7     all_users = ratings[(ratings["movieId"].isin(similar_user_recs.index)) & (ratings["rating"] > 4)]
8     all_user_recs = all_users["movieId"].value_counts() / len(all_users["userId"].unique())
9     rec_percentages = pd.concat([similar_user_recs, all_user_recs], axis=1)
10    rec_percentages.columns = ["similar", "all"]
11
12    rec_percentages["score"] = rec_percentages["similar"] / rec_percentages["all"]
13    rec_percentages = rec_percentages.sort_values("score", ascending=False)
14    return rec_percentages.head(10).merge(movies, left_index=True, right_on="movieId")[["score", "title", "genres"]]
```

```
1 import ipywidgets as widgets
2 from IPython.display import display
3
4 movie_name_input = widgets.Text(
5     value='Toy Story',
6     description='Movie Title:',
7     disabled=False
8 )
9 recommendation_list = widgets.Output()
10
11 def on_type(data):
12     with recommendation_list:
13         recommendation_list.clear_output()
14         title = data["new"]
```

```
15     if len(title) > 5:
16         results = search(title)
17         movie_id = results.iloc[0]["movieId"]
18         display(find_similar_movies(movie_id))
19
20 movie_name_input.observe(on_type, names='value')
21
22 display(movie_name_input, recommendation_list)
```

Movie Title: Toy Story

	score	title	genres
3021	18.791795	Toy Story 2 (1999)	Adventure Animation Children Comedy Fantasy
2264	8.238676	Bug's Life, A (1998)	Adventure Animation Children Comedy
2669	6.987077	Iron Giant, The (1999)	Adventure Animation Children Drama Sci-Fi
14813	6.508593	Toy Story 3 (2010)	Adventure Animation Children Comedy Fantasy IMAX
3650	6.345845	Chicken Run (2000)	Animation Children Comedy
1992	5.625199	Little Mermaid, The (1989)	Animation Children Comedy Musical Romance
2895	5.446722	Who Framed Roger Rabbit? (1988)	Adventure Animation Children Comedy Crime Fant...
3082	5.361188	Galaxy Quest (1999)	Adventure Comedy Sci-Fi
1818	5.353431	Mulan (1998)	Adventure Animation Children Comedy Drama Musi...