

Day 4}

Jan 11. 5:00pm GMT

cos\_sim = cosine\_similarity(count\_matrix)

movie-user\_likes = "Avatar"

def get\_index\_from\_title(title):

return df[df.title == title]["index"].values[0]

def get\_title\_from\_index(index):

return df[df.index == index]["title"].values[0]

movie\_index = get\_index\_from\_title(movie\_user\_likes)

similar\_movies = list(enumerate(cos\_sim[movie\_index]))

generates a tuple  
where first element is "index"  
and second is "score".

enumerate gives  
index and value

Sorted-similar-movies, sorted(similar\_movies,  
key = lambda x: x[1], reverse = True)  
 $\swarrow$   $\nwarrow$   
x[0] sim score largest to highest  
x[1]: index of movies

Print 10 similar movies

$$i = 0$$

```
i = 0
for element in sorted_similar_movies:
    print (get_title_from_index (element[0]))
```

```
print (get_title_from_index (element[0]))
```

$$i := i + 1$$

if  $i > 10$ :  
break.

for  $i$  in range(10):

```

i in range(1)
index = sorted-similar-movies[i][0]
title = get-title-from-index(index)
print(title)

```

Generate full code on paper.

1. Import libraries.

```
import pandas as pd
import numpy as np
from sklearn.feature_extraction_text import CountVectorizer
from sklearn.metrics.pairwise import cosine_similarity
```

*sklearn.feature\_extraction\_text*  
*sklearn.metrics.pairwise*

*corrections*

Helper fun

```
def get_index_from_title(title):
    return df[df.title == title]["index"].values[0]
```

```
def get_title_from_index(index):
    return df[df.index == index]["title"].values[0]
```

2. Read movie-dataset

```
df = pd.read_csv("movie-dataset.csv")
```

3. Select features of interest

```
features = ["keyword", "genre", "cast", "director"]
```

4. Create a new column in the dataset that combine all selected features

for feature in features:

df[feature]: df[feature].fillna('')

def combine\_features(row):

<sup>try:</sup>  
return row("keyword") + " " + row("genre")  
+ " " + row("cast") + " " +

row("director")

<sup>except:</sup>  
print("Error!", row)

df["combine features"] = df.apply(combine\_features,  
axis=1)

5. Create count matrix

cv = CountVectorizer()

count\_matrix = cv.fit\_transform(df["combine  
features"])

array = count\_matrix.toarray()

6. Compute the cosine-similarity b/n the movies based on count matrix

cos-sim = cosine-similarity (count-matrix)

	like	similar	count
movie			
user			

7. Define  $\text{movie\_user\_like} = \text{"Avatar"}$  <sup>movie user like</sup>

movie-index : get-index-from-title (movie-user-like)

similar\_movies : list (enumerate (cosine-sim [movie-index]))

↑  
give sim. to save

7. Sort-similar-movies : sorted (similar\_movies, key = lambda x: x[0], reverse = True)

8. List of similar movies

```
for i in range(10):
```

```
    index = sorted_similar_movies[i][0]
```

```
    title = get_title_from_index(index)
```

```
    print(title)
```