

```
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

def load_movies():
    movies = pd.read_csv('data/movies.csv')
    credits = pd.read_csv('data/credits.csv')

    # Filter movies with budget > 10M and popularity > 50
    movies = movies[movies['budget'] > 10000000 & movies['popularity'] > 50]

    # Create a dictionary mapping movie IDs to their titles
    movie_titles = {}
    for movie_id, title in movies.iterrows():
        movie_titles[movie_id] = title

    # Create a dictionary mapping movie IDs to their genres
    movie_genres = {}
    for movie_id, genres in movies.iterrows():
        movie_genres[movie_id] = genres

    # Create a dictionary mapping movie IDs to their keywords
    movie_keywords = {}
    for movie_id, keywords in movies.iterrows():
        movie_keywords[movie_id] = keywords

    # Create a dictionary mapping movie IDs to their original language
    movie_language = {}
    for movie_id, language in movies.iterrows():
        movie_language[movie_id] = language

    # Create a dictionary mapping movie IDs to their overview
    movie_overview = {}
    for movie_id, overview in movies.iterrows():
        movie_overview[movie_id] = overview

    # Create a dictionary mapping movie IDs to their production company
    movie_production = {}
    for movie_id, production in movies.iterrows():
        movie_production[movie_id] = production

    return movie_titles, movie_genres, movie_keywords, movie_language, movie_overview, movie_production

def load_credits():
    credits = pd.read_csv('data/credits.csv')

    # Create a dictionary mapping movie IDs to their cast and crew
    movie_cast_crew = {}
    for movie_id, cast_crew in credits.iterrows():
        movie_cast_crew[movie_id] = cast_crew

    return movie_cast_crew

def create_embeddings(movie_titles, movie_genres, movie_keywords, movie_language, movie_overview, movie_production, movie_cast_crew):
    # Create a TfidfVectorizer for movie titles
    vectorizer = TfidfVectorizer(stop_words='english')
    titles_embeddings = vectorizer.fit_transform(movie_titles.values).toarray()

    # Create a TfidfVectorizer for movie genres
    vectorizer = TfidfVectorizer(stop_words='english')
    genres_embeddings = vectorizer.fit_transform(movie_genres.values).toarray()

    # Create a TfidfVectorizer for movie keywords
    vectorizer = TfidfVectorizer(stop_words='english')
    keywords_embeddings = vectorizer.fit_transform(movie_keywords.values).toarray()

    # Create a TfidfVectorizer for movie language
    vectorizer = TfidfVectorizer(stop_words='english')
    language_embeddings = vectorizer.fit_transform(movie_language.values).toarray()

    # Create a TfidfVectorizer for movie overview
    vectorizer = TfidfVectorizer(stop_words='english')
    overview_embeddings = vectorizer.fit_transform(movie_overview.values).toarray()

    # Create a TfidfVectorizer for movie production
    vectorizer = TfidfVectorizer(stop_words='english')
    production_embeddings = vectorizer.fit_transform(movie_production.values).toarray()

    # Create a TfidfVectorizer for movie cast and crew
    vectorizer = TfidfVectorizer(stop_words='english')
    cast_crew_embeddings = vectorizer.fit_transform(movie_cast_crew.values).toarray()

    return titles_embeddings, genres_embeddings, keywords_embeddings, language_embeddings, overview_embeddings, production_embeddings, cast_crew_embeddings

def find_similar_movies(movie_id, titles_embeddings, genres_embeddings, keywords_embeddings, language_embeddings, overview_embeddings, production_embeddings, cast_crew_embeddings):
    # Find similar movies based on title
    titles_similarities = cosine_similarity(titles_embeddings[movie_id], titles_embeddings)
    titles_similar_movies = titles_similarities.argsort()[::-1]

    # Find similar movies based on genre
    genres_similarities = cosine_similarity(genres_embeddings[movie_id], genres_embeddings)
    genres_similar_movies = genres_similarities.argsort()[::-1]

    # Find similar movies based on keyword
    keywords_similarities = cosine_similarity(keywords_embeddings[movie_id], keywords_embeddings)
    keywords_similar_movies = keywords_similarities.argsort()[::-1]

    # Find similar movies based on language
    language_similarities = cosine_similarity(language_embeddings[movie_id], language_embeddings)
    language_similar_movies = language_similarities.argsort()[::-1]

    # Find similar movies based on overview
    overview_similarities = cosine_similarity(overview_embeddings[movie_id], overview_embeddings)
    overview_similar_movies = overview_similarities.argsort()[::-1]

    # Find similar movies based on production
    production_similarities = cosine_similarity(production_embeddings[movie_id], production_embeddings)
    production_similar_movies = production_similarities.argsort()[::-1]

    # Find similar movies based on cast and crew
    cast_crew_similarities = cosine_similarity(cast_crew_embeddings[movie_id], cast_crew_embeddings)
    cast_crew_similar_movies = cast_crew_similarities.argsort()[::-1]

    # Return the top 10 similar movies
    similar_movies = []
    for i in range(10):
        similar_movies.append({
            'movie_id': titles_similar_movies[i],
            'titles_similarity': titles_similarities[movie_id, titles_similar_movies[i]],
            'genres_similarity': genres_similarities[movie_id, genres_similar_movies[i]],
            'keywords_similarity': keywords_similarities[movie_id, keywords_similar_movies[i]],
            'language_similarity': language_similarities[movie_id, language_similar_movies[i]],
            'overview_similarity': overview_similarities[movie_id, overview_similar_movies[i]],
            'production_similarity': production_similarities[movie_id, production_similar_movies[i]],
            'cast_crew_similarity': cast_crew_similarities[movie_id, cast_crew_similar_movies[i]]
        })

    return similar_movies

# Example usage
movie_titles, movie_genres, movie_keywords, movie_language, movie_overview, movie_production, movie_cast_crew = load_movies()
movie_cast_crew = load_credits()

titles_embeddings, genres_embeddings, keywords_embeddings, language_embeddings, overview_embeddings, production_embeddings, cast_crew_embeddings = create_embeddings(movie_titles, movie_genres, movie_keywords, movie_language, movie_overview, movie_production, movie_cast_crew)

movie_id = 0
similar_movies = find_similar_movies(movie_id, titles_embeddings, genres_embeddings, keywords_embeddings, language_embeddings, overview_embeddings, production_embeddings, cast_crew_embeddings)

for i in range(10):
    print(f"Movie {i+1}: {similar_movies[i]['movie_id']} (Titles: {similar_movies[i]['titles_similarity']}, Genres: {similar_movies[i]['genres_similarity']}, Keywords: {similar_movies[i]['keywords_similarity']}, Language: {similar_movies[i]['language_similarity']}, Overview: {similar_movies[i]['overview_similarity']}, Production: {similar_movies[i]['production_similarity']}, Cast/Crew: {similar_movies[i]['cast_crew_similarity']})")
```



```
Out[109]: array([[0, 0, ..., 0, 0, 0],
                [0, 0, ..., 0, 0, 0],
                [0, 0, ..., 0, 0, 0],
                ...,
                [0, 0, ..., 0, 0, 0],
                [0, 0, ..., 0, 0, 0],
                [0, 0, ..., 0, 0, 0], dtype=int64)
```

```
In [97]: vectors[0]
```

```
Out[97]: array([0, 0, 0, ..., 0, 0, 0], dtype=int64)
```

```
In [110]: cv.get_feature_names()
```



```

'confess',
'confid',
'conflict',
'confound',
'confuse',
'confuse',
'confuse',
'connect',
'connecticut',
'connecti',
'connel',
'connote',
'conquer',
'consequa',
'consequences',
'conserve',
'conserv',
'conspiraci',
'conspiraci',
'conspiracy',
'constant',
'constantli',
'construct',
'consum',
'contact',
'contain',
'contemporari',
'content',
'contest',
'contest',
'contest',
'contract',
'contractor',
'control',
'controversi',
'convert',
'convert',
'convers',
'convince',
'convince',
'cook',
...])
```

```
In [100]: import nltk

In [101]: from nltk.stem.porter import PorterStemmer
ps = PorterStemmer()

In [103]: def stem(text):
    y = []
    for i in text.split():
        y.append(ps.stem(i))
    return " ".join(y)

In [104]: new_df['tags'] = new_df['tags'].apply(stem)

C:\python-input-106\bel8a346d89b\1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
new_df['tags'] = new_df['tags'].apply(stem)

In [104]: stem("In the 22nd century, a paraplegic Marine is dispatched to the moon Pandora on a unique mission, but becom

Out[104]: 'In the 22nd century, a paraplegic marin is dispatch to the moon pandora on a uniqu mission, but becom torn betwe
war spacecoloniz societ1 spacetravel futurist romans space alien tribe alienplanet cgi marin soldier battl lovea
ffair antiwar powerrel mindandsoul 3d samworthington zoosaldana sigourneyweav jamescameron'

In [111]: from sklearn.metrics.pairwise import cosine_similarity

In [118]: similarity = cosine_similarity(vectors)

In [119]: similarity[1]

Out[119]: array([0.08346223, 1.          , 0.06063391, ..., 0.02378257, 0.          ,
0.02615329])

In [123]: sorted(list(enumerate(similarity[0])),reverse = True,key= lambda x:x[1])[1:6]

Out[123]: [(1216, 0.2867696673382022),
(2409, 0.26901379342488517),
(3759, 0.2603130246478754),
(507, 0.255608593705383),
(539, 0.25038669783359574)]

In [ ]: def recommend(movie):
    index = new_df['title'] == movie].index[0]
    distances = sorted(list(enumerate(similarity[index])),reverse=True,key = lambda x: x[1])
    for i in distances[1:6]:
        new_df.iloc[i[0]].title

In [131]: def recommend(movie):
    movie_index = new_df['title'] == movie].index[0]
    distances = similarity[movie_index]
    movies_list = sorted(list(enumerate(distances)),reverse = True,key= lambda x:x[1])[1:6]
    for i in movies_list:
        print(new_df.iloc[i[0]].title)

In [132]: recommend('Batman Begins')

The Dark Knight
Batman
Batman
The Dark Knight Rises
10th & Wolf

In [128]: new_df.iloc[1216].title

Out[128]: 'Aliens vs Predator: Requiem'

In [ ]:
```