

آرین قره محمد زاده قشقایی

4017724001

محمد برزگر

4010724001

تمرین 5

دیتاست movie_5000_t به وسیله api کگل به کولب اضافه شده و دیتاست را دانلود کردیم در ابتدا.

```
!pip install kaggle
```

```
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Requirement already satisfied: kaggle in /usr/local/lib/python3.10/dist-packages (1.5.13)
Requirement already satisfied: six>=1.10 in /usr/local/lib/python3.10/dist-packages (from kaggle) (1.16.0)
Requirement already satisfied: certifi in /usr/local/lib/python3.10/dist-packages (from kaggle) (2022.12.7)
Requirement already satisfied: python-dateutil in /usr/local/lib/python3.10/dist-packages (from kaggle) (2.8.2)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from kaggle) (2.27.1)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from kaggle) (4.65.0)
Requirement already satisfied: python-slugify in /usr/local/lib/python3.10/dist-packages (from kaggle) (8.0.1)
Requirement already satisfied: urllib3 in /usr/local/lib/python3.10/dist-packages (from kaggle) (1.26.15)
Requirement already satisfied: text-unidecode>=1.3 in /usr/local/lib/python3.10/dist-packages (from python-slugify->kaggle) (1.3)
Requirement already satisfied: charset-normalizer<=2.0.0 in /usr/local/lib/python3.10/dist-packages (from requests->kaggle) (2.0.12)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->kaggle) (3.4)
```

```
from google.colab import files
files.upload()
```

Choose Files kaggle.json

- kaggle.json(application/json) - 66 bytes, last modified: 6/16/2023 - 100% done

Saving kaggle.json to kaggle.json

```
{'kaggle.json': b'{"username": "arianghmgh", "key": "33d467c6c0fb3830e865dd6fedeb6075"}'}
```

```
[9] !mkdir ~/.kaggle
!cp kaggle.json ~/.kaggle/
!chmod 600 ~/.kaggle/kaggle.json
```

```
[10] !kaggle datasets download -d tmdb/tmdb-movie-metadata
```

```
Downloading tmdb-movie-metadata.zip to /content
0% 0.00/8.89M [00:00<?, ?B/s]
100% 8.89M/8.89M [00:00<00:00, 123MB/s]
```

```
[11] !unzip tmdb-movie-metadata.zip
```

```
Archive: tmdb-movie-metadata.zip
  inflating: tmdb_5000_credits.csv
  inflating: tmdb_5000_movies.csv
```

```
[12] credits = pd.read_csv('tmdb_5000_credits.csv')
```

دیتاست حاوی 2 فایل `tmdb_5000_credits.csv` و `tmdb_5000_movies.csv` بود که

فایل `tmdb_5000_movies.csv` حاوی ستون های زیر بود

```
'budget', 'genres', 'homepage', 'id', 'keywords', 'original_language',  
'original_title', 'overview', 'popularity', 'production_companies',  
'production_countries', 'release_date', 'revenue', 'runtime',  
'spoken_languages', 'status', 'tagline', 'title', 'vote_average',  
'vote_count'], dtype='object'
```

از بین این ستون ها، ستون های `keywords`، `title`، `overview`، `tagline` را انتخاب و عمل پری پراکسیسینگ را انجام میدیم.

دو تابع زیر برای حذف کلمات `stem` و `stopwords` ها ساخته شده.

```
nltk.download('stopwords')  
def RemoveStopWord(clean_tokens):  
    stop_words = set(stopwords.words('english'))  
    filtered_tokens = [token for token in clean_tokens if token.lower() not in stop_words and len(token) > 1]  
    return filtered_tokens
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...  
[nltk_data] Unzipping corpora/stopwords.zip.
```

```
nltk.download('punkt')  
def stem_words(word_array):  
    stemmer = PorterStemmer()  
    stemmed_words = [stemmer.stem(word) for word in word_array]  
    return stemmed_words
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...  
[nltk_data] Unzipping tokenizers/punkt.zip.
```


از sklearn متدی برای محاسبه Tfidf به پروژه اضافه میکنیم. و به وسیله آن وکتور های کلمات را بدست می آوریم.

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer()
tfidf_matrix = vectorizer.fit_transform(Data_after_preprocessing['overview'])
feature_names = vectorizer.get_feature_names_out()
for i, doc in enumerate(Data_after_preprocessing['overview']):
    feature_index = tfidf_matrix[i, :].nonzero()[1]
    tfidf_scores = zip(feature_index, [tfidf_matrix[i, x] for x in feature_index])
    tfidf_scores = sorted(tfidf_scores, key=lambda x: (x[1]), reverse=True)

    print(f"TF-IDF scores for Document ({i}):")
    for feature_idx, score in tfidf_scores:
        print(f"{feature_names[feature_idx]}: {score}")
    print("\n")

la: 0.1080291297131482
encount: 0.106502011747466
part: 0.10511800100093662
exper: 0.10423380988896435
leav: 0.0948515367855979
endi: 0.086427717299876
soon: 0.08425484147828472
day: 0.08255471165009984
woman: 0.08155851200131921
two: 0.06848623142320492
life: 0.06164593425963316

TF-IDF scores for Document 3956:
friendsEngl: 0.2564234941204849
reduc: 0.23573413539286567
garag: 0.23573413539286567
entrepreneur: 0.22987364952813158
enabl: 0.2150467755872464
advantag: 0.20838429872251233
object: 0.20551943086261293
appar: 0.19834112092131384
invent: 0.19834112092131384
capabl: 0.19834112092131384
mess: 0.19834112092131384
devic: 0.19628179165725998
uniqu: 0.18623862185606152
highli: 0.1812529187163748
consequ: 0.17785176131559458
opportun: 0.17461370794503503
unexpect: 0.16348993378638843
accident: 0.1603480872465197
seenEngl: 0.15975400980896035
anyth: 0.15860805150960835
challeng: 0.1575060275066403
```

همین طور این متد به صورت دستی توسط ما نیز پیاده سازی شد.

```
import math
from collections import Counter

def calculate_tf(document):
    tf_scores = {}
    word_counts = Counter(document.split())
    total_words = len(document.split())
    for word, count in word_counts.items():
        tf_scores[word] = count / total_words
    return tf_scores

def calculate_idf(documents):
    idf_scores = {}
    total_documents = len(documents)
    all_words = set([word for document in documents for word in document.split()])
    for word in all_words:
        document_count = sum([1 for document in documents if word in document.split()])
        idf_scores[word] = math.log(total_documents / (1 + document_count))
    return idf_scores

def calculate_tf_idf(document, documents):
    tf_idf_scores = {}
    tf_scores = calculate_tf(document)
    idf_scores = calculate_idf(documents)
    for word in document.split():
        tf_idf_scores[word] = tf_scores[word] * idf_scores[word]
    return tf_idf_scores
```

در اینجا دو دامیومنت را به صورت رندوم انتخاب و tf idf آن را محاسبه کردیم

وکتورهای محاسبه شده را به تابع زیر دادیم

که یک معیار شباهت برحسب \cos similarity هست دادیم و مقدار شباهت این دو وکتور را بدست آوردیم.