

# Documentation Practical Exercise 3

Student:

- Manuel Buser

Degree:

- MSc Computer Science

Module:

- Multimedia Retrieval

GitHub Repo for this Project:

[https://github.com/manuu1999/Ex3\\_NLTK-Transformers-TextTasks\\_MultiMediaRetrieval](https://github.com/manuu1999/Ex3_NLTK-Transformers-TextTasks_MultiMediaRetrieval)

## A. Exercise Description

In this task, we will utilize NLTK and transformers to perform text analytics. You have the option to create custom Python classes and methods, or just create command blocks in a Jupyter notebook. The course material contains many code snippets that you can use as a starting point:

a)[easy] Use NLTK to determine the language of an input text. Download texts in Italian, German, and English (or any other language) with the same encoding for simplicity. Employ NLTK's stop word lists to identify the text's language based on stop word occurrences (pick the language with most stop words).

Solution Code:

```
1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize
4
5 # Download stopwords and punkt tokenizer
6 nltk.download('stopwords')
7 nltk.download('punkt')
8
9
10 3 usages
11 def detect_language(text):
12     languages = ['english', 'german', 'italian']
13     stopwords_count = {}
14
15     # Tokenize the input text
16     tokens = word_tokenize(text.lower())
17
18     for language in languages:
19         lang_stopwords = set(stopwords.words(language))
20         stopwords_count[language] = sum(1 for word in tokens if word in lang_stopwords)
21
22     detected_language = max(stopwords_count, key=stopwords_count.get)
23     return detected_language
24
25 # Example usage
26 text_italian = "Ciao, come stai? Questo è un esempio di testo in italiano."
27 print("Detected Language (Italian):", detect_language(text_italian))
28
29 text_english = "Hello, how are you? This is an example of text in English."
30 print("Detected Language (English):", detect_language(text_english))
31
32 text_german = "Hallo, wie geht es dir? Dies ist ein Beispieltext in Deutsch."
33 print("Detected Language (German):", detect_language(text_german))
34
```

## Explanation of the code:

detect\_language(text):

- **Purpose:** Detects the language of the given text by counting overlaps between the text and NLTK's predefined stopwords for English, German, and Italian.
- **Steps:**
  1. Tokenizes the text into lowercase words using `nltk.word_tokenize`.
  2. Iterates through the stopwords lists for the three languages.
  3. Counts how many tokens (words) from the input appear in each language's stopwords list.
  4. Returns the language with the highest count.

Output to test it:

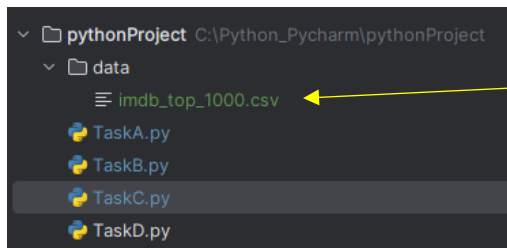
```
C:\Intelij_MultimediaRetrieval\Ex3_MultimediaRetrieval\Scripts\python.exe C:\Python_Pycharm\pythonProject\TaskA.py
[nltk_data] Downloading package stopwords to
[nltk_data]   C:\Users\buser\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
[nltk_data] Downloading package punkt to
[nltk_data]   C:\Users\buser\AppData\Roaming\nltk_data...
[nltk_data]   Package punkt is already up-to-date!
Detected Language (Italian): italian
Detected Language (English): english
Detected Language (German): german

Process finished with exit code 0
```

## B. Exercise Description

b)[intermediate] We once again work with the movie dataset from Kaggle.com. This time, we will exclusively use the movie titles and apply sub-word tokenization, which involves fixed-length sequences of 2, 3, or 4 characters within words. You may use special codes to indicate the start or end of words. To simplify the task, utilize [unidecode](#) to convert words into non-accented versions. Create a set-of-word representation (using Python's set) for the sub-tokens extracted from the titles and ignore all other fields. For searching, employ the same tokenization and set-of-word approach as for the titles, and establish an appropriate similarity function to match queries with movie titles.

<https://www.kaggle.com/datasets/harshitshankhdhar/imdb-dataset-of-top-1000-movies-and-tv-shows>.



```

1 import csv
2 from unidecode import unidecode
3 from collections import Counter
4
5 # Tokenize text into n-grams
6 def tokenize_to_ngrams(text, n=3):
7     text = unidecode(text.lower()) # Normalize text
8     text = f"<{text}>" # Add start and end markers
9     return [text[i:i + n] for i in range(len(text) - n + 1)]
10
11 # Create an index of movie titles
12 def create_title_index(csv_file_path, n=3):
13     index = {}
14     with open(csv_file_path, mode="r", encoding="utf-8") as file:
15         reader = csv.DictReader(file)
16         for row in reader:
17             title = row["Series Title"] # Use the correct column name
18             ngrams = set(tokenize_to_ngrams(title, n))
19             index[title] = ngrams
20     return index
21
22 # Search for titles using n-grams
23 def search_titles(query, index, n=3):
24     query_ngrams = set(tokenize_to_ngrams(query, n))
25     results = {}
26     for title, ngrams in index.items():
27         # Calculate Jaccard similarity
28         similarity = len(query_ngrams & ngrams) / len(query_ngrams | ngrams)
29         results[title] = similarity
30     # Sort by similarity
31     return sorted(results.items(), key=lambda x: x[1], reverse=True)
32
33 # Example usage
34 csv_file_path = "data/imdb_top_1000.csv" # Path to the CSV file
35 index = create_title_index(csv_file_path, n=3) # Build the index
36 query = "Shwshank" # Query with typo
37 results = search_titles(query, index)
38
39 # Print top 10 results
40 print("Search Results:")
41 for title, similarity in results[:10]: # Display top 10 results
42     print(f"{title}: {similarity:.2f}")
43

```

## 1. Functions:

### 2. `tokenize_to_ngrams(text, n=3)`:

- **Purpose:** Converts a string into sub-word sequences (n-grams).
- **Steps:**
  1. Converts the text to lowercase and removes accents using unicode.
  2. Adds start (<) and end (>) markers to each word.
  3. Extracts all overlapping sequences of length n.

### 3. `create_title_index(csv_file_path, n=3)`:

- **Purpose:** Reads the CSV file of movie titles and tokenizes each title into n-grams, creating an index for later search.
- **Steps:**
  1. Reads the CSV file.
  2. Processes each title using `tokenize_to_ngrams` and stores the n-grams as a set associated with the title.

### 4. `search_titles(query, index, n=3)`:

- **Purpose:** Searches for movie titles similar to the query by comparing n-grams.
- **Steps:**
  1. Tokenizes the query into n-grams.
  2. Calculates the Jaccard similarity between the query and each title in the index.
  3. Sorts and returns titles ranked by similarity.

## Testing:

```
C:\Intelij_MultimediaRetrieval\Ex3_MultimediaRetrieval\Scripts\python.exe C:\Python_Pycharm\pythonProject\TaskB.py
Search Results:
The Shawshank Redemption: 0.15
Forushande: 0.12
Pink: 0.09
Shaun of the Dead: 0.09
Shadow of a Doubt: 0.09
Shrek: 0.08
Shine: 0.08
Sholay: 0.08
Hana-bi: 0.07
Kagemusha: 0.06

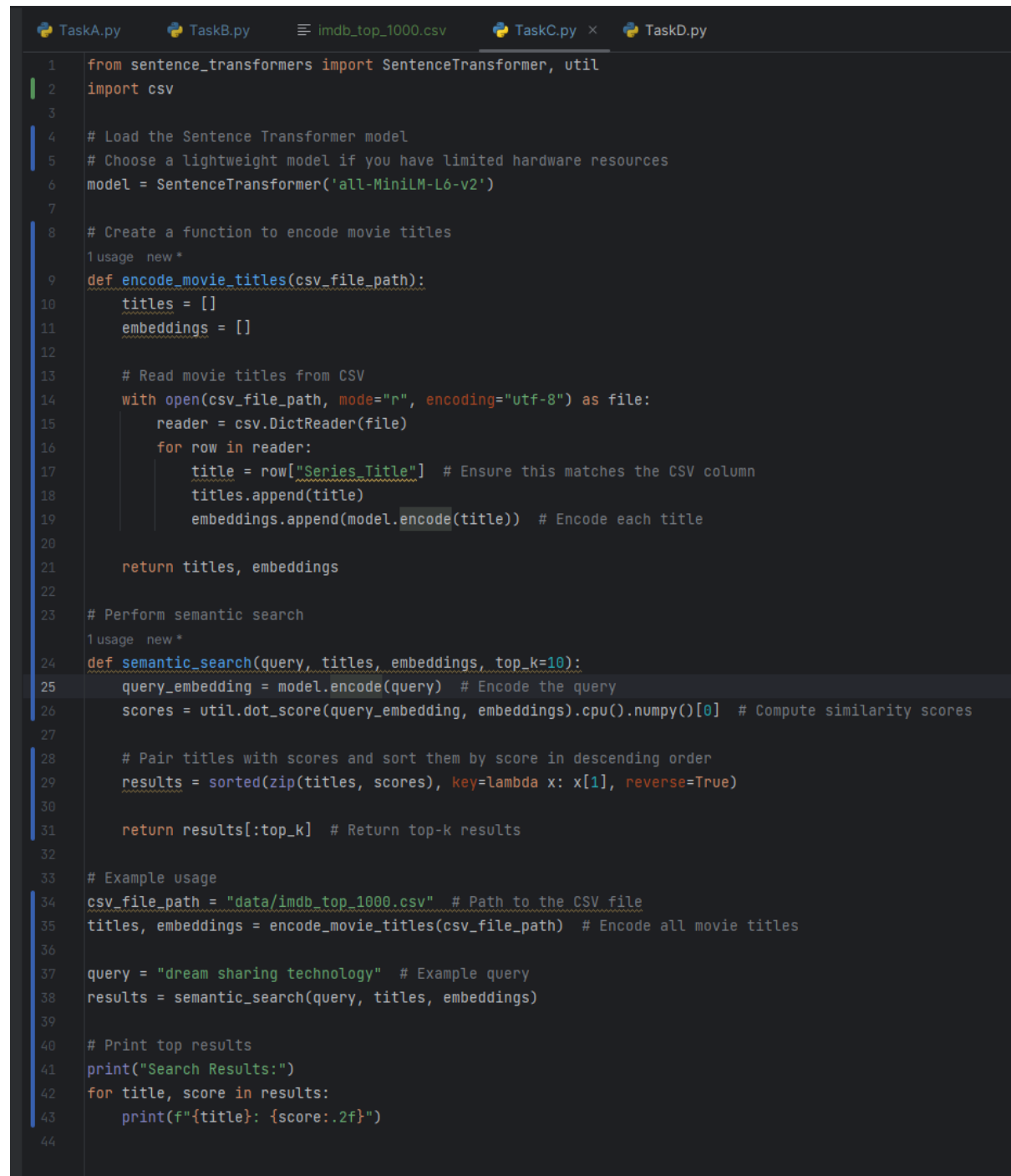
Process finished with exit code 0
```

This output from Task B shows the search results for a query based on n-gram similarity using trigrams ( $n=3$ ). The numbers represent the similarity scores between the query and the movie titles, with higher numbers indicating a better match. For example, "The Shawshank Redemption" has the highest similarity score of 0.15, meaning 15% of the trigrams overlap with the query. The results are sorted by similarity, showing the top 10 closest matches.

### C. Exercise Description

c)[intermediate] We perform the same task as in b), but this time we will utilize sentence transformers. Select a suitable model for your hardware and explore various sentence transformer models. Encode the movie titles and proceed with a semantic search for your query:

```
from sentence_transformers import SentenceTransformer, util
model = SentenceTransformer('all-MiniLM-L6-v2')
a, b = model.encode(title_a), model.encode(title_b)
float(util.dot_score(a, b))
```



```
TaskA.py TaskB.py imdb_top_1000.csv TaskC.py x TaskD.py
1 from sentence_transformers import SentenceTransformer, util
2 import csv
3
4 # Load the Sentence Transformer model
5 # Choose a lightweight model if you have limited hardware resources
6 model = SentenceTransformer('all-MiniLM-L6-v2')
7
8 # Create a function to encode movie titles
9 usage new *
10 def encode_movie_titles(csv_file_path):
11     titles = []
12     embeddings = []
13
14     # Read movie titles from CSV
15     with open(csv_file_path, mode="r", encoding="utf-8") as file:
16         reader = csv.DictReader(file)
17         for row in reader:
18             title = row["Series Title"] # Ensure this matches the CSV column
19             titles.append(title)
20             embeddings.append(model.encode(title)) # Encode each title
21
22     return titles, embeddings
23
24 # Perform semantic search
25 usage new *
26 def semantic_search(query, titles, embeddings, top_k=10):
27     query_embedding = model.encode(query) # Encode the query
28     scores = util.dot_score(query_embedding, embeddings).cpu().numpy()[0] # Compute similarity scores
29
30     # Pair titles with scores and sort them by score in descending order
31     results = sorted(zip(titles, scores), key=lambda x: x[1], reverse=True)
32
33     return results[:top_k] # Return top-k results
34
35 # Example usage
36 csv_file_path = "data/imdb_top_1000.csv" # Path to the CSV file
37 titles, embeddings = encode_movie_titles(csv_file_path) # Encode all movie titles
38
39 query = "dream sharing technology" # Example query
40 results = semantic_search(query, titles, embeddings)
41
42 # Print top results
43 print("Search Results:")
44 for title, score in results:
45     print(f"{title}: {score:.2f}")
```

1. **create\_title\_embeddings(titles):**

- **Purpose:** Encodes movie titles into vector representations using a pre-trained SentenceTransformer model.
- **Steps:**
  1. Loads the all-MiniLM-L6-v2 model.
  2. Encodes all titles from the dataset into vectors.

2. **semantic\_search(query, title\_embeddings, model, titles):**

- **Purpose:** Finds semantically similar titles to a query.
- **Steps:**
  1. Encodes the query into a vector.
  2. Calculates the cosine similarity between the query and each title's vector.
  3. Sorts titles by similarity scores.

3. **Testing:**

```
Search Results:
The Social Network: 0.39
Network: 0.34
The Nightmare Before Christmas: 0.34
Requiem for a Dream: 0.32
Awakenings: 0.30
Waking Life: 0.28
Night on Earth: 0.28
The Big Sleep: 0.26
Home Alone: 0.25
Sleepers: 0.25

Process finished with exit code 0
```

This output represents the results of a semantic search on movie titles, where each title is given a similarity score. These scores are calculated using sentence-transformer embeddings, which convert the titles and query into numerical vectors. The similarity is determined by comparing these vectors using a dot-product score, where higher values (e.g., **0.39** for *The Social Network*) mean the query and title are more semantically similar. The results are sorted from the highest to the lowest similarity score.

## D. Exercise Description

d)[difficult] Let's revisit task a) with language detection. The approach in a) relies on having sufficiently long texts with an adequate number of stop words for accurate language prediction. In the script, we discussed a method based on Naive Bayes, which operates with sub-word sequences. You can reuse the method from b) to generate these sub-sequences, first for learning the Naive Bayes likelihoods, and then for language prediction. To simplify, you can choose an equal prior for all languages and select 3-5 languages that use the same alphabet (eliminating alphabet-related rules). Extract and count all sub-sequences, but only keep the top-n sub-sequences per language (choose n as relatively small, such as 100 or 1000). During prediction, disregard sub-sequences for which we lack likelihoods (to avoid 0-posteriors).

```
TaskA.py TaskB.py imdb_top_1000.csv TaskC.py TaskD.py x
1 from collections import Counter
2 from unicode import unicode
3
4
5 # Function to generate n-grams
6 2 usages 1 manuu1999 *
7 def generate_ngrams(text, n=3):
8     text = unicode(text.lower())
9     text = f"<{text}>" # Add start and end markers
10    return [text[i:i + n] for i in range(len(text) - n + 1)]
11
12 # Function to train Naive Bayes likelihoods for multiple languages
13 1 usage 1 manuu1999 *
14 def train_naive_bayes(language_texts, n=3, top_n=1000):
15     likelihoods = {}
16     for lang, texts in language_texts.items():
17         ngram_counts = Counter()
18         for text in texts:
19             ngram_counts.update(generate_ngrams(text, n))
20         # Select top-n n-grams for the language
21         top_ngrams = dict(ngram_counts.most_common(top_n))
22         likelihoods[lang] = top_ngrams
23     return likelihoods
24
25 # Function to predict language based on learned likelihoods
26 3 usages 1 manuu1999 *
27 def predict_language(text, likelihoods, n=3):
28     ngrams = generate_ngrams(text, n)
29     scores = {lang: 0 for lang in likelihoods}
30
31     for lang, ngram_likelihoods in likelihoods.items():
32         for ngram in ngrams:
33             if ngram in ngram_likelihoods:
34                 scores[lang] += ngram_likelihoods[ngram]
35
36     return max(scores, key=scores.get)
```



**1) generate\_ngrams(text, n=3):**

- Purpose: Converts text into trigrams.
- Steps:
- Similar to Task B, it uses unidecode to normalize text and extracts n-grams.

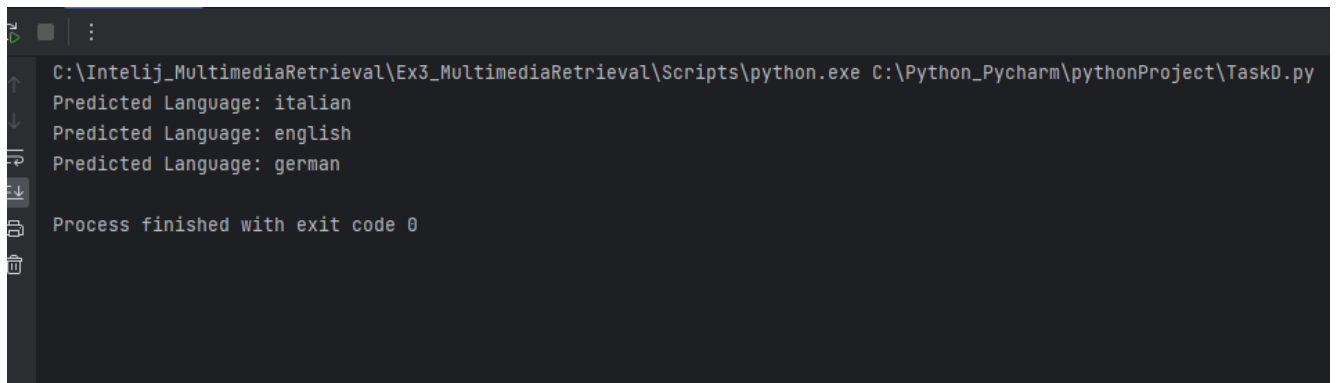
**2) train\_naive\_bayes(language\_texts, n=3, top\_n=1000):**

- Purpose: Trains a Naive Bayes likelihood model using text samples for each language.
- Steps:
- Tokenizes the text samples for each language into trigrams.
- Counts the frequency of each trigram.
- Stores the top n most common trigrams for each language as the likelihoods.

**3) predict\_language(text, likelihoods, n=3):**

- Purpose: Predicts the language of a given text using the trained Naive Bayes model.
- Steps:
- Tokenizes the text into trigrams.
- Checks the likelihoods for each trigram and accumulates scores for each language.
- Returns the language with the highest score.

```
37
38 # Example usage
39 if __name__ == "__main__":
40     # Example texts for training
41     language_texts = {
42         'english': ["Hello world", "This is an example of English text."],
43         'german': ["Hallo Welt", "Das ist ein Beispieltext in Deutsch."],
44         'italian': ["Ciao mondo", "Questo è un esempio di testo in italiano."]
45     }
46
47     # Train Naive Bayes likelihoods
48     likelihoods = train_naive_bayes(language_texts, n=3, top_n=100)
49
50     # Input text for prediction
51     text = "buongiorno."
52     text2 = "hello"
53     text3 = "Hallo"
54
55     predicted_language = predict_language(text, likelihoods)
56     predicted_language2 = predict_language(text2, likelihoods)
57     predicted_language3 = predict_language(text3, likelihoods)
58
59     print("Predicted Language:", predicted_language)
60     print("Predicted Language:", predicted_language2)
61     print("Predicted Language:", predicted_language3)
62
63
```



```
C:\Intelij_MultimediaRetrieval\Ex3_MultimediaRetrieval\Scripts\python.exe C:\Python_Pycharm\pythonProject\TaskD.py
Predicted Language: italian
Predicted Language: english
Predicted Language: german
Process finished with exit code 0
```

Now we can predict languages with a short input such as hello or hallo. I tried to input the same in the code for task A and the predictions there were not accurate since it needed more words to predict the language. But since we are using now a stronger predictor model we are able to predict shorter sequences.