

1. Cover Page

- Course title: Natural Language Processing (Spring 2025)
- topic name: **Assignment 1 + Assignment 2**
- student's name: Nurzhan Mussabekov
- date: 14.02.2025

2. Introduction

- Brief overview of NLP and Deep Learning:

Natural Language Processing (NLP) is a specialized field within artificial intelligence that focuses on enabling computers to understand, interpret, and generate human language. NLP encompasses a variety of tasks, including text analysis, language translation, sentiment analysis, and named entity recognition. In recent years, the integration of deep learning techniques has significantly advanced NLP. Deep learning models, particularly those based on neural networks such as transformers, can capture complex patterns and contextual nuances in language data, leading to improvements in accuracy and efficiency for many NLP applications.

- Importance of Python libraries for NLP (NLTK, spaCy, transformers):

NLTK provides a comprehensive set of tools for text processing, including tokenization, lemmatization, and stopword removal. It is especially useful for educational purposes and prototyping traditional NLP techniques.

spaCy offers robust capabilities for tokenization, lemmatization, part-of-speech tagging, and named entity recognition (NER). Its context-aware models and built-in visualization tools (like displacy) make it well-suited for production-level NLP tasks.

Transformers library has revolutionized NLP by providing access to state-of-the-art pretrained transformer models such as BERT, GPT, and RoBERTa. These models excel in capturing contextual information and semantic meaning in text, enabling advanced applications like text vectorization and sentiment analysis with deep learning techniques.

3. Implementation and Code Snippets

- Ex1. Text Preprocessing with NLTK and spaCy:

task 1. Tokenize a sample paragraph using NLTK and spaCy:

A brief explanation of the task: The goal of this experiment is to preprocess a sample paragraph using two popular NLP libraries, **NLTK** and **spaCy**. The preprocessing steps include tokenization, lemmatization, and stopwords removal. The outputs from both libraries are compared to analyze their differences.

Code snippets with appropriate comments:

Code #1:

```
import requests
from bs4 import BeautifulSoup

url =
"https://en.wikipedia.org/wiki/Natural_language_processing"

response = requests.get(url)
soup = BeautifulSoup(response.text, "html.parser")

text = soup.find_all("p")[0].get_text()
print("Sample Text:", text)
```

Code #2:

```
import nltk #nltk library

nltk.download('punkt_tab') #pretrained model for tokenizing text.

from nltk.tokenize import word_tokenize #Import the word_tokenize function

nltk_tokens = word_tokenize(text) #tokenize the text

print("NLTK Tokenization:", nltk_tokens)
```

Code #3:

```
import spacy #import spaCy

nlp = spacy.load("en_core_web_sm") #load the small English model

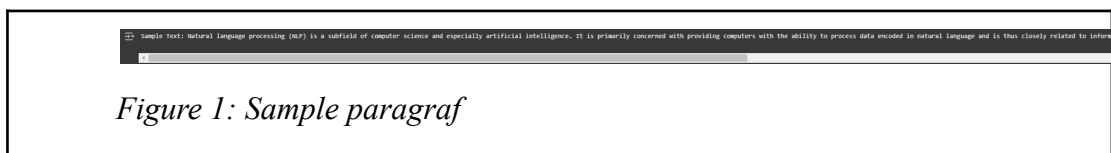
doc = nlp(text) #process the text

spacy_tokens = [token.text for token in doc] #get tokens from the processed text

print("spaCy Tokenization:", spacy_tokens)
```

Images of the executed code:

output of Code #1:



output of Code #2:

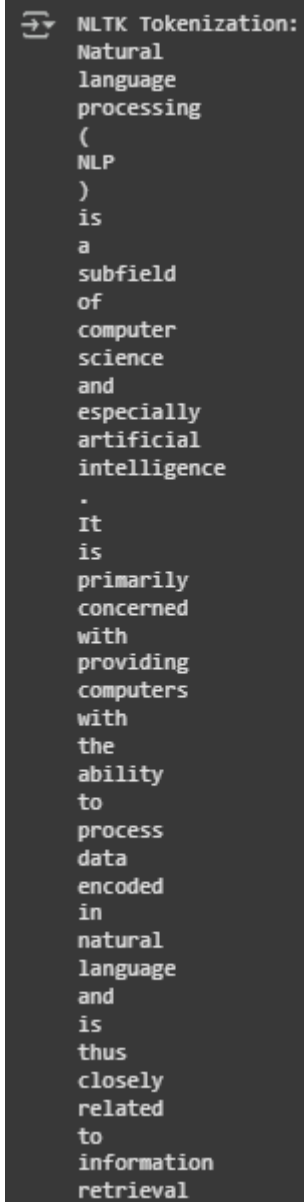
 NLTK Tokenization:
Natural
language
processing
(
NLP
)
is
a
subfield
of
computer
science
and
especially
artificial
intelligence
.
It
is
primarily
concerned
with
providing
computers
with
the
ability
to
process
data
encoded
in
natural
language
and
is
thus
closely
related
to
information
retrieval

Figure 2: Tokenization using NLTK

Output of Code #3:

```
→ /usr/local/lib/python3.11/dist-  
warnings.warn(Warnings.W111)  
spaCy Tokenization:  
Natural  
language  
processing  
(  
NLP  
)  
is  
a  
subfield  
of  
computer  
science  
and  
especially  
artificial  
intelligence  
.  
It  
is  
primarily  
concerned  
with  
providing  
computers  
with  
the  
ability  
to  
process  
data  
encoded  
in  
natural  
language  
and  
is  
thus  
closely  
related  
to  
information  
retrieval
```

Figure 3: Tokenization using SpaCy

task 2. Perform lemmatization and stopwords removal using both libraries:

A brief explanation of the task: In this task, we preprocess a sample text by applying lemmatization and stopwords removal using two popular NLP libraries, NLTK and spaCy. Lemmatization reduces words to their base or dictionary form (e.g., "running" becomes "run"), which helps in normalizing the text, while stopwords removal eliminates common words (like "the", "and", "is") that often do not add meaningful context to the analysis. By performing these steps with both NLTK and spaCy, we can

compare their approaches and results to understand their differences in handling linguistic nuances.

Code snippets with appropriate comments:

Code #1:

```
nltk.download('wordnet') # download WordNet corpus for lemmatization
nltk.download('omw-1.4') #additional language data for WordNet
nltk.download('stopwords') #download stopwords list

from nltk.tokenize import word_tokenize # import tokenization function
from nltk.corpus import stopwords #import stopwords
from nltk.stem import WordNetLemmatizer # import WordNet lemmatizer

nltk_tokens = word_tokenize(text) #tokenize the text

lemmatizer = WordNetLemmatizer() #initialize the lemmatizer
nltk_stopwords = set(stopwords.words('english')) #create a set of english
stopwords

nltk_processed = [lemmatizer.lemmatize(token) for token in nltk_tokens if
token.lower() not in nltk_stopwords] # lemmatize tokens and remove stopwords

print("NLTK Processed:", nltk_processed)
```

Code #2:


```
doc = nlp(text) # process the input text

spacy_processed = [token.lemma_ for token in doc if not token.is_stop]
#lemmatize tokens and remove stopwords using spacy

print("spaCy Processed:", spacy_processed)
```

Images of the executed code:

Output of Code #1:



```
NLTK Processed:
Natural
language
processing
(
NLP
)
subfield
computer
science
especially
artificial
intelligence
.
primarily
concerned
providing
computer
ability
process
data
encoded
natural
language
thus
closely
related
information
retrieval
,
knowledge
representation
computational
linguistics
,
subfield
linguistics
.
Typically
data
collected
text
corpus
,
using
```

Figure 4: lemmatization and stopword removal using NLTK library.

Output of Code #2:


 spaCy Processed:
natural
language
processing
(
NLP
)
subfield
computer
science
especially
artificial
intelligence
.
primarily
concern
provide
computer
ability
process
datum
encode
natural
language
closely
relate
information
retrieval
,
knowledge
representation
computational
linguistic
,
subfield
linguistic
.
typically
datum
collect
text
corpora
,
rule
-
base
,
statistical
neural
-
base
approach
machine
learning
deep

Figure 5: lemmatization and stopword removal using SpaCy library.

Ex 2. Named Entity Recognition (NER) with spaCy

task 1. Use spaCy's pre-trained NER model to extract named entities from a given text.

A brief explanation of the task: This task uses spaCy's pre-trained English model to process a given text and extract named entities (e.g., persons, organizations, locations). The model identifies and labels the entities automatically based on its built-in NER pipeline.

Code snippets with appropriate comments:

Code #1:

```
import spacy #import the spaCy library

nlp = spacy.load("en_core_web_sm") #Load spaCy's small English model

text = "Almaty,[a] formerly Alma-Ata,[b] is the largest city in Kazakhstan, with a population exceeding two million residents within its metropolitan area.[8] Located in the foothills of the Trans-Ili Alatau mountains in southern Kazakhstan, near the border with Kyrgyzstan, Almaty stands as a pivotal center of culture, commerce, finance and innovation. The city is nestled at an elevation of 700–900 metres (2,300–3,000 feet), with the Big Almaty and Small Almaty rivers running through it, originating from the surrounding mountains and flowing into the plains. Almaty is the second-largest city in Central Asia and the third-largest in the Eurasian Economic Union (EEU)."
```

```
doc = nlp(text) #process the text

print("Named Entities:")
for ent in doc.ents:
    print(ent.text, "->", ent.label_)
```

Images of the executed code:

Output of Code #1:

```

/usr/local/lib/python3.11/dist-packag
warnings.warn(Warnings.W111)
Named Entities:
Almaty,[a -> CARDINAL
Kazakhstan -> GPE
two million -> CARDINAL
the Trans-Ili Alatau -> ORG
Kazakhstan -> GPE
Kyrgyzstan -> GPE
Almaty -> GPE
700-900 metres -> QUANTITY
2,300-3,000 feet -> QUANTITY
the Big Almaty -> ORG
Small Almaty -> PERSON
Almaty -> ORG
second -> ORDINAL
Central Asia -> LOC
third -> ORDINAL
the Eurasian Economic Union -> ORG

```

Figure 6: NER Extraction using spaCy

task 2. Visualize the named entities using displacy.

A brief explanation of the task: This task utilizes spaCy's built-in visualization tool, displacy, to render the named entities in the text. The visualization highlights the entities with color-coded labels, making it easier to see how the model has categorized parts of the text.

Code snippets with appropriate comments:

```

from spacy import displacy # Import displacy for visualization

displacy.render(doc, style="ent", options={"distance": 120}) #visualize the named
entities in the processed document using displacy

```

Images of the executed code:

Almaty,[a CARDINAL] formerly Alma-Ata,[b] is the largest city in Kazakhstan GPE, with a population exceeding two million CARDINAL residents within its metropolitan area.[8] Located in the foothills of the Trans-Ili Alatau ORG mountains in southern Kazakhstan GPE, near the border with Kyrgyzstan GPE, Almaty GPE stands as a pivotal center of culture, commerce, finance and innovation. The city is nestled at an elevation of 700-900 metres QUANTITY (2,300-3,000 feet QUANTITY), with the Big Almaty ORG and Small Almaty PERSON rivers running through it, originating from the surrounding mountains and flowing into the plains. Almaty ORG is the second ORDINAL - largest city in Central Asia LOC and the third ORDINAL -largest in the Eurasian Economic Union ORG (EEU).

Figure 7: Visualization of Named Entities using displacy

Ex 3. Text Vectorization using Transformers:

task 1. Load a pretrained transformer model from Hugging Face (e.g., bert-base-uncased):

A brief explanation of the task: This task loads the pretrained BERT model (bert-base-uncased) along with its corresponding tokenizer from Hugging Face's Transformers library. The model is used for later steps in text vectorization.

Code snippets with appropriate comments:

Code #1:

```
from transformers import AutoTokenizer, AutoModel #import necessary classes from
Hugging Face

tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased") #load the tokenizer
for bert-base-uncased

model = AutoModel.from_pretrained("bert-base-uncased") #load the BERT model for
bert-base-uncased
```

Images of the executed code:

Output of Code #1:

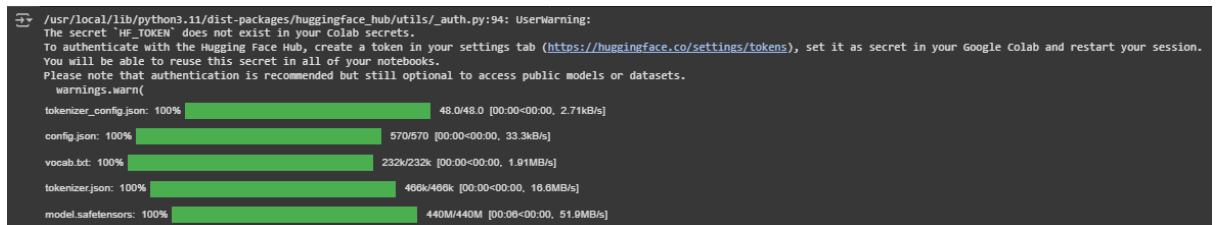


Figure 8: Loading Transformer Model

task 2. Tokenize and encode a sample sentence using the tokenizer:

A brief explanation of the task: This task tokenizes and encodes a sample sentence using the loaded tokenizer. The encoding converts the text into numerical representations (token IDs and attention masks) suitable for the transformer model.

Code snippets with appropriate comments:

Code #1:

Code snippets with appropriate comments:

Code #1:

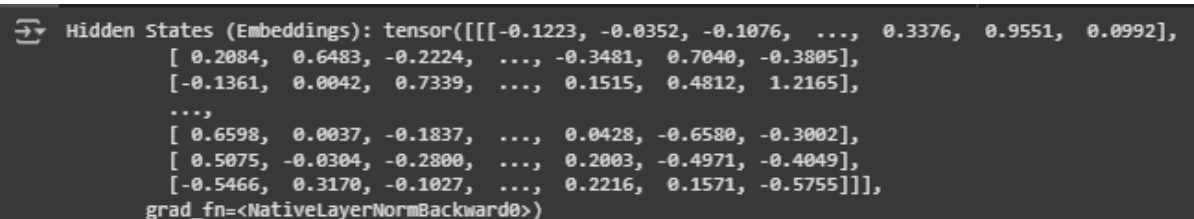
```
outputs = model(**encoded_input) # pass the encoded input through the model

hidden_states = outputs.last_hidden_state # The last hidden state contains the word
embeddings for each token

print("Hidden States (Embeddings):", hidden_states)
```

Images of the executed code:

Output of Code #1:



```
Hidden States (Embeddings): tensor([[[[-0.1223, -0.0352, -0.1076, ..., 0.3376, 0.9551, 0.0992],
[ 0.2084, 0.6483, -0.2224, ..., -0.3481, 0.7040, -0.3805],
[-0.1361, 0.0042, 0.7339, ..., 0.1515, 0.4812, 1.2165],
...,
[ 0.6598, 0.0037, -0.1837, ..., 0.0428, -0.6580, -0.3002],
[ 0.5075, -0.0304, -0.2800, ..., 0.2003, -0.4971, -0.4049],
[-0.5466, 0.3170, -0.1027, ..., 0.2216, 0.1571, -0.5755]]]],
grad_fn=<NativeLayerNormBackward0>)
```

Figure 10: Extracted Word Embeddings

Ex 4. Sentiment Analysis with Transformers:

task 1. Use the pipeline module from Hugging Face to perform sentiment analysis on different sentences:

A brief explanation of the task: This task uses the Hugging Face pipeline module to perform sentiment analysis on multiple sentences. The pipeline loads a pretrained sentiment analysis model (by default, a model fine-tuned on sentiment tasks) and applies it to input text, returning a sentiment label (e.g., "POSITIVE" or "NEGATIVE") along with a confidence score.

Code snippets with appropriate comments:

Code #1:

```
import pandas as pd

from transformers import pipeline # Import Hugging Face's pipeline for sentiment
analysis
```

```
url =
"https://raw.githubusercontent.com/SK7here/Movie-Review-Sentiment-Analysis/refs/
heads/master/IMDB-Dataset.csv"

df = pd.read_csv(url) #load the dataset directly from the URL

reviews = df["review"].head(10).tolist() #for demonstration selected the first 10
movie reviews from the dataset

sentiment_pipeline = pipeline("sentiment-analysis") #initialize the sentiment analysis
pipeline with a pretrained model

print("Sentiment Analysis using Transformers on Movie Reviews:")
for review in reviews:
    result = sentiment_pipeline(review)
    print("Review:", review)
    print("Result:", result, "\n")
```

Images of the executed code:

Output of Code #1:

```

No model was supplied, defaulted to distilbert/distilbert-base-unc
Using a pipeline without specifying a model name and revision in p
Device set to use cpu
Sentiment Analysis using Transformers on Movie Reviews:
Review: One of the other reviewers has mentioned that after watchi
Result: [{'label': 'NEGATIVE', 'score': 0.5121204853057861}]

Review: A wonderful little production. <br /><br />The filming tec
Result: [{'label': 'POSITIVE', 'score': 0.9993765950202942}]

Review: I thought this was a wonderful way to spend time on a too
Result: [{'label': 'POSITIVE', 'score': 0.9991896748542786}]

Review: Basically there's a family where a little boy (Jake) think
Result: [{'label': 'NEGATIVE', 'score': 0.9991722106933594}]

Review: Petter Mattei's "Love in the Time of Money" is a visually
Result: [{'label': 'POSITIVE', 'score': 0.9988238215446472}]

Review: Probably my all-time favorite movie, a story of selflessne
Result: [{'label': 'POSITIVE', 'score': 0.9997684359550476}]

Review: I sure would like to see a resurrection of a up dated Seah
Result: [{'label': 'POSITIVE', 'score': 0.9160271286964417}]

Review: This show was an amazing, fresh & innovative idea in the 7
Result: [{'label': 'NEGATIVE', 'score': 0.9996374845504761}]

Review: Encouraged by the positive comments about this film on her
Result: [{'label': 'NEGATIVE', 'score': 0.9997221827507019}]

Review: If you like original gut wrenching laughter you will like
Result: [{'label': 'POSITIVE', 'score': 0.9997014403343201}]

```

Figure 11: Sentiment Analysis with Transformers

task 2. Compare results with traditional text-processing approaches:

A brief explanation of the task: In this task, we implement a traditional sentiment analysis approach using a Naive Bayes classifier. We vectorize the same 10 movie reviews using a bag-of-words representation (via CountVectorizer) and then train a MultinomialNB classifier. Finally, we compare the predictions from this classical method with the results from the transformer-based approach. rule-based model that provides sentiment scores (positive, neutral, negative, and compound). The comparison highlights the differences between modern model-based sentiment analysis and traditional methods.

Code snippets with appropriate comments:

Code #1:

```

import pandas as pd
from sklearn.feature_extraction.text import CountVectorizer

```



```

from sklearn.naive_bayes import MultinomialNB

url =
"https://raw.githubusercontent.com/SK7here/Movie-Review-Sentiment-Analysis/refs/
heads/master/IMDB-Dataset.csv"

df = pd.read_csv(url) #load the dataset directly from the URL

reviews = df["review"].head(10).tolist() #for demonstration selected the first 10
movie reviews from the dataset

sentiments = df["sentiment"].head(10).tolist() #sentiment

# Vectorize the text data using CountVectorizer (bag-of-words model)
vectorizer = CountVectorizer(stop_words="english")
X = vectorizer.fit_transform(reviews)

#initialize and train the Naive Bayes classifier
nb_classifier = MultinomialNB()
nb_classifier.fit(X, sentiments)

predictions = nb_classifier.predict(X) #predict sentiments using the trained classifier

print("Naive Bayes Sentiment Analysis:")
for review, actual, predicted in zip(reviews, sentiments, predictions):
    print("Review:", review)
    print("Actual Sentiment:", actual, "-> Predicted Sentiment:", predicted)
    print("-" * 80)

```

Images of the executed code:

Output of Code #1:

```

Naive Bayes Sentiment Analysis:
Review: One of the other reviewers has mentioned that after watching just 1 Oz episode
Actual Sentiment: positive -> Predicted Sentiment: positive
-----
Review: A wonderful little production. <br /><br />The filming technique is very u
Actual Sentiment: positive -> Predicted Sentiment: positive
-----
Review: I thought this was a wonderful way to spend time on a too hot summer week
Actual Sentiment: positive -> Predicted Sentiment: positive
-----
Review: Basically there's a family where a little boy (Jake) thinks there's a zomb
Actual Sentiment: negative -> Predicted Sentiment: negative
-----
Review: Petter Mattei's "Love in the Time of Money" is a visually stunning film to
Actual Sentiment: positive -> Predicted Sentiment: positive
-----
Review: Probably my all-time favorite movie, a story of selflessness, sacrifice an
Actual Sentiment: positive -> Predicted Sentiment: positive
-----
Review: I sure would like to see a resurrection of a up dated Seahunt series with
Actual Sentiment: positive -> Predicted Sentiment: positive
-----
Review: This show was an amazing, fresh & innovative idea in the 70's when it fir
Actual Sentiment: negative -> Predicted Sentiment: negative
-----
Review: Encouraged by the positive comments about this film on here I was looking
Actual Sentiment: negative -> Predicted Sentiment: negative
-----
Review: If you like original gut wrenching laughter you will like this movie. If y
Actual Sentiment: positive -> Predicted Sentiment: positive
-----

```

Figure 12: Sentiment Analysis using Naive Bayes.

4. Results and Discussion

- Images of results with proper explanations:
- Figure 1: Figure 1: Sample paragraph - shows output of code to scrape the sample paragraph from wikipedias NLP page by using BeautifulSoup
- *Figure 2: Tokenization using NLTK* - shows code to tokenize the input text into a list of word tokens using NLTK's pre-trained 'punkt_tab' model.
- *Figure 3: Tokenization using SpaCy* - shows code to tokenize the input text into a list of word tokens using spaCy's pre-trained English language model:
- *Figure 4: lemmatization and stopwords removal using NLTK library* - shows code to lemmatization and stopwords removal using NLTK library:
- *Figure 5: lemmatization and stopwords removal using SpaCy library.* - shows code to lemmatization and stopwords removal using Spacy library:
- Figure 6: This image displays output of the spaCy model identifying various named entities. The output lists each entity followed by its label.

Figure 7: This image shows the visual representation of the named entities extracted from the sample text. The displacy visualization clearly marks entities using different colors and labels. This graphical view helps in quickly identifying and understanding the types of entities detected by the model.

Figure 8: This image shows that the BERT model and its tokenizer have been successfully loaded from Hugging Face, which is essential for tokenizing and vectorizing the text in subsequent tasks.

Figure 9: This image demonstrates the conversion of the sample text into a dictionary of tensors. The keys (e.g., `input_ids` and `attention_mask`) represent the numerical format of the text, which is required for processing by the transformer model.

Figure 10: The screenshot illustrates the output tensor representing the hidden states of the model. Each vector in this tensor corresponds to a token from the input text, providing a dense representation that captures its semantic meaning.

Figure 11: This image shows that the Hugging Face transformer model has analyzed each sentence and returned sentiment results for the IMDB review dataset. For example, a positive sentence displays like `{'label': 'POSITIVE', 'score': 0.99}`, This confirms the model's ability to gauge sentiment based on its learned patterns.

Figure 12: The screenshot shows the results from the Naive Bayes classifier. For each movie review, the output displays the original text, the actual sentiment from the

dataset, and the sentiment predicted by the classifier. This demonstrates a more traditional text processing approach to sentiment analysis.

- Comparisons between different approaches:

Ex 1 task 3. Compare the outputs and explain the differences: Comparison of Figure 2 and Figure 3 show that there is no difference when word tokenization whenever use Spacy or NLTK. Comparison of Figure 3 and Figure 4 show that *spaCy* uses contextual information from its full NLP pipeline to determine the correct lemma while *NLTK*'s lemmatizer works in a more isolated manner and may default to less accurate results.

Ex 4 task 2. Sentiment Analysis with Transformers: Comparison of the Transformer model in Figure 11 and the Traditional Naive Bayes model in Figure 12 shows that while the transformer-based method generally offers higher accuracy, this comparison highlights the differences in performance and approach between deep learning and classical machine learning techniques.

5. Conclusion

- Summary of key findings:

Text Preprocessing with NLTK and spaCy: Both libraries successfully tokenized the sample paragraph, but subtle differences were observed. NLTK's tokenizer (using `punkt_tab`) generally produced tokens based on predefined rules, whereas spaCy's tokenizer was more context-aware and often handled punctuation and compound words differently. In terms of lemmatization and stopword removal, spaCy's integrated pipeline automatically provided context-sensitive lemmatization and used a refined stopword list. In contrast, NLTK's WordNetLemmatizer required additional handling (such as proper POS tagging) to achieve comparable results.

Named Entity Recognition (NER) with spaCy: spaCy's pre-trained NER model effectively extracted named entities from the given text, identifying a variety of entity types (e.g., organizations, locations, and persons). The use of displacy for visualization provided an intuitive and immediate understanding of the entities detected, demonstrating the strength of spaCy's integrated visualization capabilities.

Text Vectorization using Transformers: Loading a pretrained transformer model (such as `bert-base-uncased`) from Hugging Face allowed for a modern approach to text vectorization. Tokenizing and encoding a sample sentence produced numerical representations (token IDs and attention masks) which were then used to extract rich word embeddings from the model's hidden states. These embeddings capture deeper semantic and contextual information compared to traditional bag-of-words representations, illustrating the advantage of transformer-based vectorization.

Sentiment Analysis with Transformers: Utilizing Hugging Face's pipeline for sentiment analysis enabled state-of-the-art sentiment classification directly on movie reviews, with the model providing sentiment labels along with confidence scores. When compared with a traditional Naive Bayes classifier that employed bag-of-words vectorization, the transformer-based approach was more robust and effective at capturing nuances in the sentiment of complex text.

- Insights on which methods were more effective and why:

Modern Transformer-Based Methods: Transformer models (such as BERT) consistently provided more effective results, particularly in tasks like text vectorization and sentiment analysis. Their ability to capture contextual nuances and generate rich, dense embeddings makes them superior for complex tasks where deep semantic understanding is required. The end-to-end pipelines provided by Hugging Face also simplify the process and reduce the need for extensive manual preprocessing.

spaCy for Preprocessing and NER: spaCy's integrated NLP pipeline (including tokenization, lemmatization, and NER) proved to be efficient and accurate. Its context-aware features and built-in visualization tools (displacy) facilitate rapid prototyping and offer production-ready performance. In contrast, while NLTK is useful for educational purposes and offers a modular approach to NLP tasks, it requires additional steps and manual intervention to match the performance of spaCy in real-world applications.

Traditional Approaches: Traditional methods like Naive Bayes classifiers are simpler and computationally less expensive. However, they generally fall short when compared to transformer-based models in capturing the intricacies of language, especially for tasks involving nuanced sentiment or complex syntactic structures.

\

6. References

- <https://spacy.io/api/entityrecognizer>
- NLTK Official Documentation: <https://www.nltk.org/>
- BERT Model (bert-base-uncased) Documentation:
<https://huggingface.co/bert-base-uncased>
- Scikit-learn Documentation (for methods like Naive Bayes and CountVectorizer):
<https://scikit-learn.org/stable/>