| Subject : (Code) | Natural Language Processing (Group B) (IS355TBC) | Semester: 5th BE | |
|---|---|---|---|
| Date : | Duration: 120 minutes | Staff: Dr.Sindhu D V | |
| Name : | USN : | Section : | A |

512

| Q. No | Part-A | Marks | Level | CO |
|---|---|---|---|---|
| 1 | True | 1 | 2 | 2 |
| 2 | Transformers model the textual context but not in a sequential manner. Given a word in the input, it prefers to look at all the words around it and represent each word with respect to its context. | 1 | 1 | 1 |
| 3 | NLTK | 2 | 3 | 2 |
| 4 | Because of aspects like ambiguity, the need for contextual information, and idioms | 1 | 1 | 1 |
| 5 | The number of times a word appears in a document | 2 | 3 | 1 |
| 6 | [A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Za-z]{2,} | 2 | 2 | 1 |
| 7 | ☐ **Tokenization**: <br> • Tokenization is the process of **splitting text** into smaller units called **tokens**. These tokens could be words, characters, or subwords, depending on the level of tokenization. The main goal is to break down a sentence or a piece of text into manageable chunks for further analysis. <br> • **Example**: <br>    o Input text: "I love NLP." <br>    o Tokens: ['I', 'love', 'NLP', '.'] <br> ☐ **Lemmatization**: <br> • Lemmatization is the process of reducing words to their **base or root form**, known as a **lemma**. Unlike stemming (which simply removes suffixes), lemmatization considers the context of the word and applies proper linguistic rules to transform words to their base form. <br> • **Example**: <br>    o Input word: "running" <br>    o Lemmatized form: "run" <br>    o Input word: "better" <br>    o Lemmatized form: "good" | 2 | 1 | 1 |

| Q. No | Answer all questions | Marks | Level | CO |
|-------|---------------------|-------|-------|-----|
| 1a | Despite such tremendous success, DL is still not the silver bullet for all NLP tasks when it comes to the industrial applications. Some of the key reasons for this are as follows:<br>Overfitting on small datasets<br>Few-shot learning and synthetic data generation<br>Domain adaption<br>Interpretable models<br>Common sense and world knowledge. | 5 | 2 | 3 |
| 1b | <br>1. Speech recognition<br>2. Natural Language Understanding<br>3. Dialog management<br>4. Response generation<br>5. Speech synthesis | 5 | 2 | 1 |
| 2 | Data is a heart of any ML system. In most of the industrial projects, it is often the data becomes the bottleneck.<br>▪ Use public dataset<br>▪ Scrape data<br>▪ Product intervention<br>▪ Data augmentation<br>1. Synonym replacement<br>2. Back translation<br>3. Bigram flipping<br>4. Replacing entities<br>5. Adding noise to data | 10 | 4 | 2 |
| 3 | 1. Machine Translation | 2m× | 2 | 1 |

| | | | | |
|---|---|---|---|---|
| | 2. Sentimental analysis | 5=10 | | |
| 4a | import nltk<br>from nltk.corpus import stopwords<br>from nltk.tokenize import word_tokenize<br>import string<br><br># Ensure necessary resources are downloaded<br>nltk.download('punkt')<br>nltk.download('stopwords')<br><br># Define the function<br>def load_and_clean_corpus(file_path):<br>   # Load the corpus from a file<br>   with open(file_path, 'r') as file:<br>     text = file.read()<br><br>   # Tokenize the text into words<br>   tokens = word_tokenize(text)<br><br>   # Convert to lowercase<br>   tokens = [word.lower() for word in tokens]<br><br>   # Remove punctuation<br>   tokens = [word for word in tokens if word.isalnum()]<br><br>   # Remove stopwords<br>   stop_words = set(stopwords.words('english'))<br>   cleaned_tokens = [word for word in tokens if word not in stop_words]<br><br>   return cleaned_tokens<br><br># Test the function<br>file_path = "sample_corpus.txt"  # Replace with your file's path<br>cleaned_corpus = load_and_clean_corpus(file_path)<br><br>print("Cleaned Tokens:", cleaned_corpus) | 6 | 1,2 | [1] |
| 4b | There are two approaches<br>   1. A classical NLP<br>   2. DL pipeline | 4 | 3 | 3 |

| 5a. | import re | 10 | 3 | 2 |
|---|---|---|---|---|

```
import re
import nltk
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer

# Preprocessing function
def preprocess_review(review):
    # Remove non-alphanumeric characters
    review = re.sub(r'\W', ' ', review)

    # Convert to lowercase
    review = review.lower()

    # Tokenization
    tokens = nltk.word_tokenize(review)

    # Remove stopwords
    stop_words = set(stopwords.words('english'))
    tokens = [word for word in tokens if word not in stop_words]

    # Lemmatization
    lemmatizer = WordNetLemmatizer()
    tokens = [lemmatizer.lemmatize(word) for word in tokens]

    return ' '.join(tokens)

# Preprocess all reviews
processed_reviews = [preprocess_review(review) for review in reviews]
```

```
from sklearn.feature_extraction.text import CountVectorizer

# Initialize CountVectorizer
vectorizer = CountVectorizer()

# Transform reviews into a Bag of Words representation
X = vectorizer.fit_transform(processed_reviews)
from sklearn.model_selection import train_test_split
from sklearn.naive_bayes import MultinomialNB
from sklearn.metrics import classification_report

# Sample labels (1 for positive, 0 for negative)
# In a real-world scenario, these would be obtained through manual labeling or
pre-labeled data.
labels = [1 if 'good' in review or 'excellent' in review else 0 for review in
processed_reviews]

# Split the data into train and test sets
X_train, X_test, y_train, y_test = train_test_split(X, labels, test_size=0.2,
random_state=42)

# Train Naive Bayes classifier
classifier = MultinomialNB()
classifier.fit(X_train, y_train)

# Test the classifier
y_pred = classifier.predict(X_test)

# Evaluate the model
print(classification_report(y_test, y_pred))
```

| Marks Distribution | Particulars | | CO1 | CO2 | CO3 | CO4 | L1 | L2 | L3 | L4 | L5 | L6 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Test | Max Marks | 28 | 19 | 13 | | 6 | 26 | 18 | 10 | | |

BT-Blooms Taxonomy, CO-Course Outcomes, M-Marks