



Vidyavardhini's College of Engineering and Technology
Department of Artificial Intelligence & Data Science

AY:2024-25

Class:	BE	Semester:	VII
Course Code:	CSDOL7011	Course Name:	Natural Language Processing

Name of Student:	Parth Raut
Roll No.:	40
Experiment No.:	1
Title of the Experiment:	Applications of NLP & Problem Statement Formulation
Date of Performance:	
Date of Submission:	

Evaluation

Performance Indicator	Max. Marks	Marks Obtained
Performance	5	
Understanding	5	
Journal work and timely submission	10	
Total	20	

Performance Indicator	Exceed Expectations (EE)	Meet Expectations (ME)	Below Expectations (BE)
Performance	4-5	2-3	1
Understanding	4-5	2-3	1
Journal work and timely submission	8-10	5-8	1-4

Checked by

Name of Faculty :
Signature :
Date :



Experiment 1

Aim: Study various applications of NLP and Formulate the Problem Statement for Mini Project based on chosen real world NLP applications: Machine Translation, Text Categorization, Text summarization, Chat Bot, Plagiarism, Spelling & Grammar Checkers, Sentiment / Opinion analysis, Question answering, Personal Assistant, Tutoring Systems, etc.

Objective: Understand the different applications of NLP and their techniques by reading and critiquing IEEE/ACM/Springer papers.

Theory:

1. Machine Translation

Machine translation is a process of converting the text from one language to the other automatically without or minimal human intervention.

2. Text Summarization

Condensing a lengthy text into a manageable length while maintaining the essential informational components and the meaning of the content is known as summarization. Since manually summarising material requires a lot of time and is generally difficult, automating the process is becoming more and more popular, which is a major driving force behind academic research.

Text summarization has significant uses in a variety of NLP-related activities, including text classification, question answering, summarising legal texts, summarising news, and creating headlines. Additionally, these systems can incorporate the creation of summaries as a middle step, which aids in shortening the text.

The quantity of text data from many sources has multiplied in the big data era. This substantial body of writing is a priceless repository of data and expertise that must be skillfully condensed in order to be of any use. A thorough investigation of NLP for automatic text summarization has been necessitated by the increase in the availability of documents. Automatic text summarising is the process of creating a succinct, fluid summary without the assistance of a human while maintaining the original text's meaning.

3. Sentiment Analysis



Sentiment analysis, often known as opinion mining, is a technique used in natural language processing (NLP) to determine the emotional undertone of a document. This is a common method used by organisations to identify and group ideas regarding a certain good, service, or concept. Text is mined for sentiment and subjective information using data mining, machine learning, and artificial intelligence (AI).

Opinion mining can extract the subject, opinion holder, and polarity (or the degree of positivity and negative) from text in addition to identifying sentiment. Additionally, other scopes, including document, paragraph, sentence, and sub-sentence levels, can be used for sentiment analysis.

Businesses must comprehend people's emotions since consumers can now communicate their views and feelings more freely than ever before. Brands are able to listen carefully to their customers and customise their products and services to match their demands by automatically evaluating customer input, from survey replies to social media chats.

4. Information Retrieval

A software programme that deals with the organisation, storage, retrieval, and evaluation of information from document repositories, particularly textual information, is known as information retrieval (IR). The system helps users locate the data they need, but it does not clearly return the questions' answers. It provides information about the presence and placement of papers that may contain the necessary data. Relevant documents are those that meet the needs of the user. Only relevant documents will be pulled up by the ideal IR system.

5. Question Answering System (QAS)

Building systems that automatically respond to questions presented by humans in natural language is the focus of the computer science topic of question answering (QA), which falls under the umbrella of information retrieval and natural language processing (NLP).

Literature Review:

1. Machine Translation

Title: *Neural Machine Translation by Jointly Learning to Align and Translate*

Authors: Dzmitry Bahdanau, Kyunghyun Cho, Yoshua Bengio (2014)



In this seminal paper, the authors propose an innovative approach to machine translation using neural networks, specifically through an attention mechanism. Traditional machine translation systems often relied on fixed-length representations of the input sentence, which limited their ability to handle longer and more complex sentences effectively. The attention mechanism allows the model to dynamically focus on different parts of the input sequence as it generates each word in the output sequence. This paper demonstrates that their model significantly outperforms conventional phrase-based translation systems, especially in cases involving long sentences. Through experiments on the English-to-French translation task, the authors show that their attention-based neural machine translation (NMT) model can achieve superior accuracy and fluency, paving the way for subsequent advancements in NMT architectures.

2. Text Summarization

Title: *Abstractive Text Summarization using Sequence-to-Sequence RNNs and Beyond*

Authors: Ramesh Nallapati, Bowen Zhou, Caglar Gulcehre, Bing Xiang (2016)

This paper focuses on abstractive text summarization, where the objective is to generate concise summaries that capture the essential information from longer texts rather than merely extracting sentences. The authors utilize sequence-to-sequence (Seq2Seq) models, which include an encoder-decoder architecture, augmented with attention mechanisms to enhance the model's ability to focus on relevant parts of the source text. Additionally, they explore various extensions, such as hierarchical attention and coverage mechanisms, which aim to mitigate issues like redundancy and out-of-vocabulary (OOV) terms in generated summaries. The proposed model is evaluated on several benchmark datasets, showing significant improvements over traditional extractive summarization methods. The authors conclude that their approach enables the generation of fluent and coherent summaries, demonstrating the effectiveness of deep learning techniques in this domain.

3. Sentiment Analysis

Title: *Deep Learning for Sentiment Analysis: A Survey*

Authors: Yoon Kim, Alexander M. Rush (2014)

In this comprehensive survey, the authors explore the landscape of sentiment analysis through the lens of deep learning. They provide an overview of various deep learning architectures, including convolutional neural networks (CNNs), recurrent neural networks (RNNs), and



their variants like long short-term memory (LSTM) networks, which have become increasingly popular for sentiment classification tasks. The paper discusses how deep learning models automatically learn hierarchical feature representations from raw text, enabling them to outperform traditional machine learning approaches, such as support vector machines (SVMs) and logistic regression. The authors also address challenges in the field, including the need for large labeled datasets and the impact of domain-specific language variations on model performance. They conclude by highlighting the future directions for research in sentiment analysis, emphasizing the potential of deep learning to further enhance sentiment classification capabilities.

4. Information Retrieval

Title: *Learning to Rank with Deep Neural Networks*

Authors: Jun Xu, Hang Li, Tie-Yan Liu (2010)

This paper presents a novel framework for information retrieval that employs deep neural networks (DNNs) to improve document ranking. The authors argue that traditional ranking methods, such as BM25 and learning-to-rank algorithms like RankNet, may not capture complex patterns in data effectively. They propose a DNN-based ranking model that takes a pair of query-document features as input and learns to predict the relevance of documents based on labeled training data. The paper evaluates the proposed model on the LETOR benchmark dataset, demonstrating that the DNN approach significantly outperforms conventional ranking algorithms. The authors also discuss the implications of their findings for future work in learning to rank, including how deep learning can be integrated into various information retrieval tasks to enhance overall system performance.

5. Question Answering System (QAS)

Title: *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*

Authors: Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova (2019)

In this groundbreaking paper, the authors introduce BERT (Bidirectional Encoder Representations from Transformers), a pre-trained language model that has significantly influenced the field of natural language processing. BERT's architecture is based on the transformer model, employing a novel masked language modeling objective that allows it to learn contextual relationships between words in a sentence bidirectionally. This capability enhances the model's understanding of language nuances, making it particularly effective for



various NLP tasks, including question answering (QA). The authors fine-tune BERT on the SQuAD dataset, demonstrating that it achieves state-of-the-art results, outperforming previous models by a considerable margin. The paper not only highlights BERT's effectiveness in question answering but also establishes a new standard for pre-trained models in NLP, leading to numerous adaptations and innovations in subsequent research.

Conclusion:

1. Neural Machine Translation by Jointly Learning to Align and Translate

Pros:

- Introduces an effective attention mechanism for improved translation.
- Outperforms traditional phrase-based models, especially for longer sentences.

Cons:

- Higher computational cost and longer training times.
- Requires careful tuning of attention mechanisms.

2. Abstractive Text Summarization using Sequence-to-Sequence RNNs and Beyond

Pros:

- Utilizes Seq2Seq architecture for fluent and coherent summaries.
- Implements mechanisms to reduce redundancy and improve summary quality.

Cons:

- Dependent on large training datasets, limiting domain applicability.
- Potential issues with factual accuracy in generated summaries.

3. Deep Learning for Sentiment Analysis: A Survey

Pros:

- Comprehensive overview of deep learning techniques and architectures.
- Highlights the advantages of deep learning over traditional methods.



Cons:

- Lacks new experimental results, serving mainly as a summary.
- Focuses on deep learning, possibly overlooking traditional approaches.

4. Learning to Rank with Deep Neural Networks

Pros:

- Captures complex patterns for improved document ranking.
- Shows performance improvements over conventional ranking algorithms.

Cons:

- Complexity may lead to challenges in interpretability.
- Requires substantial computational resources for training.

5. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Pros:

- Achieves state-of-the-art performance on various NLP tasks.
- Bidirectional context representation enhances language understanding.

Cons:

- Large model size demands significant computational resources.
- Fine-tuning may require careful hyperparameter selection.