

Assignment 6

Neural Network



Session: 2021 – 2025

Submitted by:

Muhammad Umair Shahid

2021-CS-144

Supervised by:

Dr. Khaldoon Syed Khurshid

Department of Computer Science

University of Engineering and Technology Lahore

Pakistan

Contents

1	Introduction	3
2	System Architecture	3
2.1	Components	3
3	Implementation Details	4
3.1	Libraries Used	4
3.2	Code Snippets and Explanations	4
3.2.1	Noun Detection and Tokenization	4
3.2.2	Bag-of-Words Vectorization	5
3.2.3	Neural Network Implementation	6
4	User Interface Design	8
4.1	Features	8
4.2	Code Snippet for GUI Components	8
5	Real-World Scenarios and Applications	10
5.1	Academic Research	10
5.2	Legal Document Analysis	10
5.3	Enterprise Knowledge Management	10
6	Discussion and Potential Enhancements	11
6.1	Neural Network Optimization	11
6.2	Performance Improvements	11
6.3	Enhanced NLP Techniques	11
6.4	User Interface Enhancements	12
7	Conclusion	12

Abstract

The Document Query System is an innovative application designed to facilitate efficient retrieval of information from a vast repository of text documents. Leveraging natural language processing (NLP) techniques for noun extraction, Bag-of-Words (BoW) vectorization, and a simple neural network, the system provides users with a user-friendly graphical interface to perform queries and obtain relevant documents swiftly. This report delves into the system's architecture, implementation details, real-world applications, and potential enhancements to bolster its effectiveness and scalability.

1 Introduction

In the age of information overload, efficiently retrieving relevant documents from extensive datasets is paramount. Traditional keyword-based search systems often fall short in understanding the contextual relevance of documents. The Document Query System aims to bridge this gap by integrating NLP techniques with machine learning algorithms to deliver more accurate and context-aware search results. This system is particularly beneficial in domains such as academic research, legal document analysis, and enterprise knowledge management, where precise information retrieval is crucial.

2 System Architecture

The system architecture comprises several interconnected components:

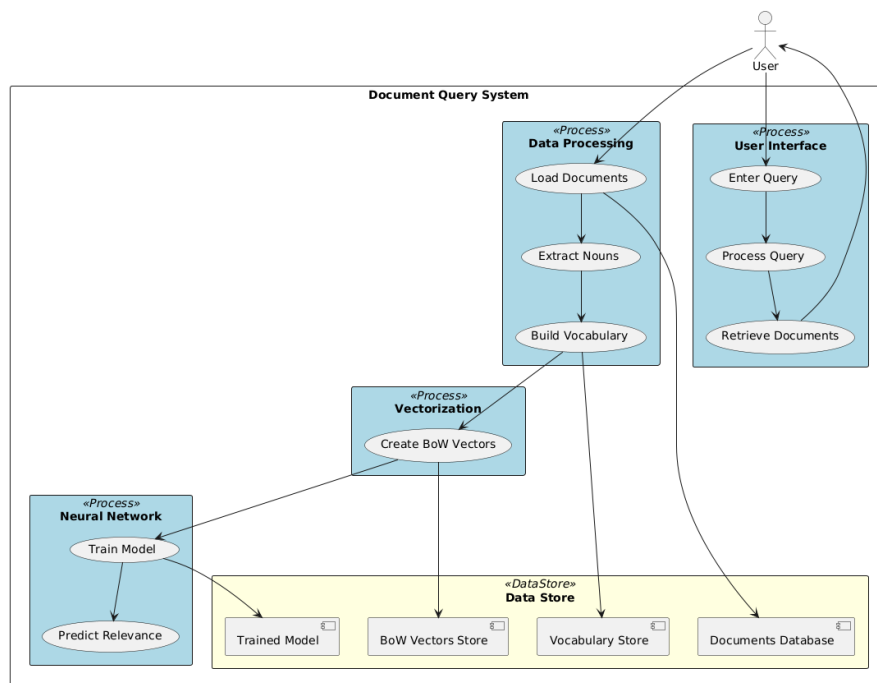


FIGURE 1: System Architecture Overview

2.1 Components

- **Data Processing Module:** Handles the loading and preprocessing of text documents, including noun extraction and vocabulary building.
- **Vectorization Module:** Converts processed text into numerical representations using the Bag-of-Words model.
- **Neural Network Module:** Implements a simple neural network for training and predicting relevance scores.

- **User Interface:** Provides a graphical interface for users to input queries and view results.

3 Implementation Details

The system is implemented in Python, utilizing several libraries to streamline development and enhance functionality.

3.1 Libraries Used

- **os:** Facilitates interaction with the operating system, enabling file operations such as reading directories and files.
- **random:** Generates random numbers, essential for initializing neural network weights and creating random labels during training.
- **re:** Utilizes regular expressions for text preprocessing, including cleaning and tokenization.
- **tkinter:** Constructs the graphical user interface (GUI), providing widgets for user interaction.
- **math.exp:** Although imported, it remains unused in the current implementation and can be omitted for efficiency.

3.2 Code Snippets and Explanations

3.2.1 Noun Detection and Tokenization

Accurate noun extraction is pivotal for building a meaningful vocabulary. The system employs a heuristic approach to identify nouns based on predefined criteria.

```
1 def is_noun(word, previous_word=None):
2     domain_proper_nouns = {
3         "ai", "ml", "nlp", "iot", "blockchain",
4         ↪ "cybersecurity",
5         "dataframe", "algorithm", "tensorflow", "pytorch",
6         "quantum", "neural", "cloud", "microservice", "api",
7         "database", "software", "hardware", "network"
8     }
9     trivial_words = {
10         "and", "in", "of", "to", "for", "the", "a", "an",
```

```
11         "from", "by", "as", "however", "when", "moreover"
12     }
13
14     noun_suffixes = [
15         "tion", "ics", "ism", "logy", "ment",
16         "ance", "ence", "ship", "ity", "ness",
17         "ware", "er", "or", "ist", "able"
18     ]
19
20     clean_word = word.strip(",.!?\\'").lower()
21
22     if clean_word in domain_proper_nouns:
23         return True
24
25     if clean_word in trivial_words:
26         return False
27
28     if word[0].isupper() and clean_word not in trivial_words:
29         return True
30
31     if any(clean_word.endswith(suffix) for suffix in
32            ↪ noun_suffixes):
33         return True
34
35     return False
```

Explanation: The 'is_noun' function determines whether a given word qualifies as a noun.

It checks against a set of domain-specific proper nouns, filters out trivial words, assesses capitalization, and evaluates noun suffixes to make this determination.

3.2.2 Bag-of-Words Vectorization

Converting text into numerical vectors is essential for machine learning algorithms to process and analyze data effectively.

```
1 def create_bow_vector(document, vocabulary):
2     tokens = document.split()
3     vector = [0] * len(vocabulary)
4     for token in tokens:
5         if token in vocabulary:
```

```

6         index = vocabulary.index(token)
7         vector[index] += 1
8     return vector

```

Explanation: The ‘create_bow_vector’ function generates a BoW vector for a given document by counting the frequency of each vocabulary term within the document. This vectorization facilitates the representation of documents in a format suitable for machine learning models.

3.2.3 Neural Network Implementation

A simple two-layer neural network is implemented to train on the BoW vectors and predict relevance scores.

```

1 class NeuralNetwork:
2     def __init__(self, input_dim, hidden_dim, output_dim):
3         self.W1 = [[random.uniform(-0.01, 0.01)
4                     for _ in range(hidden_dim)] for _ in
5                     range(input_dim)]
6         self.B1 = [0.0 for _ in range(hidden_dim)]
7         self.W2 = [[random.uniform(-0.01, 0.01) for _ in
8                     range(output_dim)]
9                     for _ in range(hidden_dim)]
10        self.B2 = [0.0 for _ in range(output_dim)]
11
12    def relu(self, x):
13        return [max(0, xi) for xi in x]
14
15    def relu_derivative(self, x):
16        return [1 if xi > 0 else 0 for xi in x]
17
18    def forward(self, x):
19        self.Z1 = [sum(x[i] * self.W1[i][j] for i in
20                      range(len(x))
21                      ) + self.B1[j] for j in
22                      range(len(self.B1))]
23        self.A1 = self.relu(self.Z1)
24        self.Z2 = [sum(self.A1[j] * self.W2[j][k]
25                      for j in range(len(self.A1))
26                      ) + self.B2[k] for k in
27                      range(len(self.B2))]

```

```

22         return self.Z2
23
24     def backward(self, x, y_true, y_pred, learning_rate):
25         dZ2 = y_pred - y_true
26         dW2 = [dZ2 * a1 for a1 in self.A1]
27         dB2 = dZ2
28
29         dA1 = [dZ2 * self.W2[j][0] for j in
30                 ↪ range(len(self.A1))]
31         dZ1 = [dA1[j] * self.relu_derivative([self.Z1[j]])[0]
32                 ↪ for j in range(len(dA1))]
33
34         dW1 = [[dZ1[j] * x[i] for j in range(len(dZ1))] for i
35                 ↪ in range(len(x))]
36         dB1 = dZ1
37
38         for j in range(len(self.W2)):
39             for k in range(len(self.W2[0])):
40                 self.W2[j][k] -= learning_rate * dW2[j]
41
42         for k in range(len(self.B2)):
43             self.B2[k] -= learning_rate * dB2
44
45         for i in range(len(self.W1)):
46             for j in range(len(self.W1[0])):
47                 self.W1[i][j] -= learning_rate * dW1[i][j]
48
49         for j in range(len(self.B1)):
50             self.B1[j] -= learning_rate * dB1[j]

```

Explanation: The 'NeuralNetwork' class encapsulates a simple neural network with one hidden layer. It includes methods for the forward pass, activation functions (ReLU), and the backward pass for weight updates using gradient descent.

4 User Interface Design

A user-friendly interface is crucial for the system's usability. The GUI, built using Tkinter, provides an intuitive platform for users to input queries and view results seamlessly.

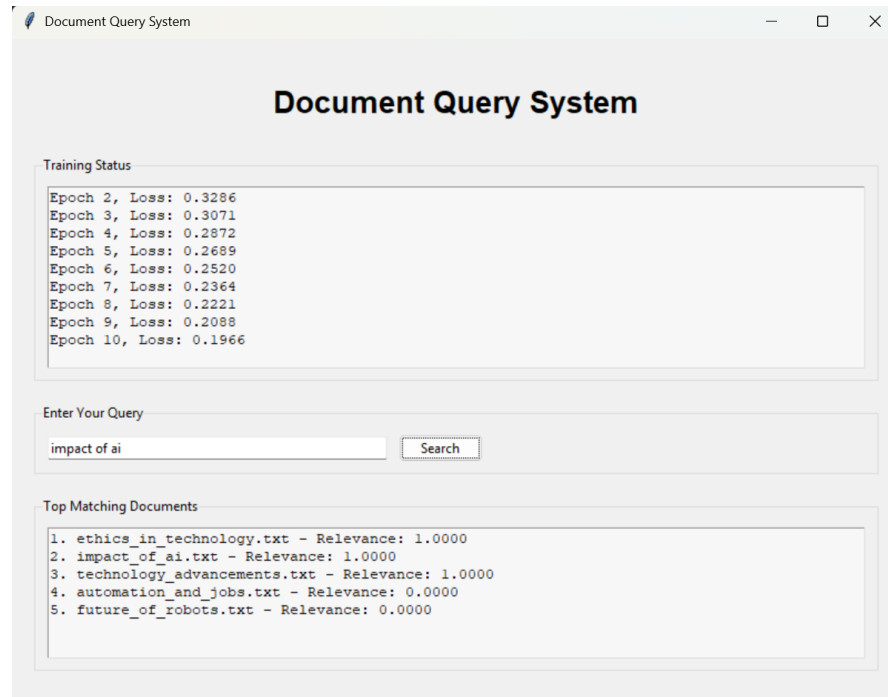


FIGURE 2: Document Query System GUI

4.1 Features

- **Query Input:** Users can enter their search queries in a dedicated input field.
- **Search Button:** Initiates the search operation based on the entered query.
- **Training Status Display:** Shows real-time updates on the training progress of the neural network.
- **Results Display:** Presents the top matching documents along with their relevance scores.
- **Exit Button:** Allows users to gracefully exit the application.

4.2 Code Snippet for GUI Components

```
1 def create_widgets(self):
2     main_frame = ttk.Frame(self.root, padding="20 20 20 20")
3     main_frame.pack(fill=tk.BOTH, expand=True)
4
```

```
5     title_label = ttk.Label(
6         main_frame, text="Document Query System",
7         ↪ font=("Helvetica", 20, "bold"))
8     title_label.pack(pady=20)
9
10    training_frame = ttk.LabelFrame(
11        main_frame, text="Training Status", padding="10 10 10
12        ↪ 10")
13    training_frame.pack(fill=tk.BOTH, expand=False, pady=10)
14    self.training_box = tk.Text(
15        training_frame, height=10, width=70, state='normal',
16        ↪ bg='#f7f7f7', wrap="word")
17    self.training_box.pack(fill=tk.BOTH, expand=True)
18    self.training_box.delete(1.0, tk.END)
19    self.training_box.config(state='disabled')
20
21    query_frame = ttk.LabelFrame(
22        main_frame, text="Enter Your Query", padding="10 10 10
23        ↪ 10")
24    query_frame.pack(fill=tk.X, expand=False, pady=10)
25
26    self.query_entry = ttk.Entry(query_frame, width=50)
27    self.query_entry.pack(side=tk.LEFT, padx=(0, 10))
28    self.query_entry.delete(0, tk.END) # Ensure it's clear
29    self.query_entry.configure(state='normal')
30    self.query_entry.focus_set()
31
32    search_button = ttk.Button(
33        query_frame, text="Search", command=self.search_query)
34    search_button.pack(side=tk.LEFT)
35
36    results_frame = ttk.LabelFrame(
37        main_frame, text="Top Matching Documents", padding="10
38        ↪ 10 10 10")
39    results_frame.pack(fill=tk.BOTH, expand=True, pady=10)
40    self.result_box = tk.Text(
41        results_frame, height=10, width=70, state='normal',
42        ↪ bg='#f7f7f7', wrap="word")
43    self.result_box.pack(fill=tk.BOTH, expand=True)
44    self.result_box.delete(1.0, tk.END)
```

```
39     self.result_box.config(state='disabled')
40
41     exit_frame = ttk.Frame(main_frame, padding="10 0 0 0")
42     exit_frame.pack(fill=tk.X, expand=False)
43     exit_button = ttk.Button(
44         exit_frame, text="Exit", command=self.root.quit)
45     exit_button.pack(side=tk.RIGHT, pady=(10, 0))
```

Explanation: The ‘create_widgets’ method constructs the GUI components, organizing them into frames and embedding functionality such as query input, search initiation, result display, and application exit.

5 Real-World Scenarios and Applications

The Document Query System is versatile and can be tailored to various real-world applications:

5.1 Academic Research

Researchers can utilize the system to sift through vast libraries of academic papers, extracting relevant studies based on specific queries. For instance, a researcher investigating “quantum algorithms” can quickly locate pertinent documents, enhancing the efficiency of literature reviews.

5.2 Legal Document Analysis

Law firms dealing with extensive legal documents can employ the system to retrieve case laws, statutes, and legal opinions relevant to particular legal queries, thereby streamlining case preparation processes.

5.3 Enterprise Knowledge Management

Organizations with large repositories of internal documents, reports, and manuals can implement the system to facilitate knowledge retrieval, ensuring that employees can access necessary information promptly.

6 Discussion and Potential Enhancements

While the current implementation lays a solid foundation, several enhancements can be made to improve the system's performance, accuracy, and user experience.

6.1 Neural Network Optimization

The existing neural network uses random labels, which hinders its ability to learn meaningful patterns. To enhance its effectiveness:

- **Meaningful Labels:** Integrate labels that reflect actual relevance scores or document categories.
- **Advanced Architectures:** Explore more sophisticated neural network architectures or employ pretrained models like BERT for better semantic understanding.
- **Training Enhancements:** Implement techniques such as learning rate scheduling, dropout, and regularization to improve training stability and prevent overfitting.

6.2 Performance Improvements

Optimizing the system for scalability and efficiency is crucial, especially when dealing with large datasets.

- **Vector Operations:** Utilize libraries like NumPy for vectorized operations to accelerate computations.
- **Efficient Data Structures:** Replace list-based vocabulary lookups with dictionaries to achieve constant-time complexity.
- **Parallel Processing:** Implement multithreading or multiprocessing to handle data loading and model training concurrently, enhancing responsiveness.

6.3 Enhanced NLP Techniques

Incorporating advanced NLP methodologies can significantly improve the system's accuracy in understanding and processing text.

- **Part-of-Speech Tagging:** Employ libraries like spaCy or NLTK for precise noun extraction and grammatical analysis.
- **Word Embeddings:** Move beyond BoW by integrating word embeddings (e.g., Word2Vec, GloVe) to capture semantic relationships between words.

- **TF-IDF Weighting:** Incorporate Term Frequency-Inverse Document Frequency (TF-IDF) to weigh the importance of terms more effectively.

6.4 User Interface Enhancements

Improving the GUI can lead to a more intuitive and efficient user experience.

- **Dynamic Dataset Selection:** Allow users to select the dataset directory through a file dialog, enhancing flexibility.
- **Progress Indicators:** Integrate progress bars or real-time status updates to inform users about ongoing operations like training.
- **Interactive Results:** Enable features such as clicking on results to open documents or view excerpts directly within the application.

7 Conclusion

The Document Query System presents a robust framework for efficient information retrieval from extensive text corpora. By integrating NLP techniques with machine learning algorithms and a user-friendly GUI, the system offers a seamless experience for users across various domains. While the current implementation serves as a foundational model, targeted enhancements in neural network training, performance optimization, advanced NLP integration, and user interface design can elevate the system's effectiveness and scalability, making it a valuable tool in the realm of information retrieval.