

Classification of Documents Using Graph-Based Features and KNN



Session: 2021 – 2025

Submitted by:

Muhammad Abdullah 2021-CS-82

Supervised by:

Sir Waqas Ali

Department of Computer Science
University of Engineering and Technology Lahore
Pakistan

Introduction

Background

Document classification is a fundamental task in natural language processing (NLP) and information retrieval. Its primary goal is to categorize documents into predefined topics or classes based on their content. Traditionally, vector-based models such as term frequency-inverse document frequency (TF-IDF) and bag-of-words (BoW) have been widely used for document classification. These models represent documents as high-dimensional vectors, where each dimension corresponds to a unique term, and classification is performed based on the similarity between these vectors.

Limitations of Vector-Based Models

While vector-based models have proven effective in many cases, they may not fully capture the semantic relationships between terms within documents. For instance, two documents may contain similar terms but express different concepts, leading to misclassification. Additionally, vector-based models suffer from the curse of dimensionality, especially when dealing with large vocabularies, which can impact classification accuracy and computational efficiency.

Motivation for Graph-Based Approach

In recent years, there has been a growing interest in graph-based methods for document classification. Graphs provide a more flexible and expressive representation of document content, allowing us to capture not only the presence of terms but also their relationships and interactions. By modeling documents as graphs, we can leverage graph algorithms and techniques to extract meaningful features and improve classification accuracy.

Significance of the Project

This project builds upon concepts from foundational papers on graph-based document classification and maximal common subgraphs (MCS) for graph comparison. By investigating graph structures to capture the inherent relationships between terms within documents, I aim to enhance classification accuracy beyond traditional vector-based models. Our approach offers hands-on experience with graph theory and machine learning, fostering skills in data representation, algorithm implementation, and analytical thinking in the context of document classification.

Objectives

The primary objective of this project is to develop a system that can classify documents into predefined topics by representing each document as a directed graph, identifying common subgraphs, and applying the K-Nearest Neighbors (KNN) algorithm based on graph similarity measures. Through this project, I seek to explore the following key components:

- Data collection and preparation
- Graph construction
- Feature extraction via common subgraphs
- Classification with KNN
- Evaluation of classification performance

Methodology

Data Collection and Preparation

Data Collection

We collected data from various online sources related to three assigned topics: travel, fashion, and diseases. Each topic consists of 15 pages of text, with approximately 300 words per page. The data collection process involved web scraping using the `requests` library to fetch HTML content from URLs and then parsing the HTML using `BeautifulSoup` to extract the text.

```
1  import requests
2  from bs4 import BeautifulSoup
3
4
5  def scrape(url_base):
6      response = requests.get(url_base)
7      if response.status_code == 200:
8          soup = BeautifulSoup(response.content, 'html.parser')
9          paragraphs = soup.find_all('p')
10         text = '\n'.join([p.text for p in paragraphs])
11         return text
12     else:
13         print(f"Failed to fetch data for Status code: {response.status_code}")
14         return None
15
16
17 def scrape_and_save_data(urls):
18     for i, url in enumerate(urls):
19         raw_text = scrape(url)
20         if raw_text:
21             word_count = len(raw_text.split())
22             file_name = f"data_{i}.txt"
23             with open(file_name, 'w', encoding='utf-8') as file:
24                 file.write(raw_text)
25                 print(f"Scraped data from {url} and saved to {file_name} and word count is {word_count}")
26
```

Data Preparation

After obtaining the raw text data, I performed preprocessing steps to prepare the data for further analysis. The preprocessing steps included:

- **Tokenization:** Breaking the text into individual words or tokens using the `nltk.tokenize.word_tokenize` function.
- **Stopword Removal:** Removing common stopwords (e.g., "the", "is", "and") using the `nltk.corpus.stopwords` module.
- **Stemming:** Reducing words to their root form using the Porter stemming algorithm implemented in the `nltk.stem.PorterStemmer` class.

```

1 import nltk
2 from nltk.corpus import stopwords
3 from nltk.tokenize import word_tokenize
4 from nltk.stem import PorterStemmer
5 import string
6 import os
7
8 # Download NLTK resources (run once)
9 # nltk.download('punkt')
10 # nltk.download('stopwords')
11
12 # Initialize stemmer and stopwords
13 stemmer = PorterStemmer()
14 stop_words = set(stopwords.words('english'))
15
16
17 def preprocess_text(raw_text):
18     # Tokenization
19     tokens = word_tokenize(raw_text.lower())
20     # Removing punctuations and stopwords, and stemming
21     processed_tokens = []
22     for token in tokens:
23         if token not in string.punctuation and token not in stop_words:
24             processed_token = stemmer.stem(token)
25             processed_tokens.append(processed_token)
26
27     # Remove duplicate words
28     unique_tokens = list(set(processed_tokens))
29
30     return unique_tokens
31
32
33 def preprocess_files(folder):
34     for file_name in os.listdir(folder):
35         if file_name.endswith('.txt'):
36             with open(os.path.join(folder, file_name), 'r', encoding='utf-8') as file:
37                 raw_text = file.read()
38                 preprocessed_text = preprocess_text(raw_text)
39                 with open(os.path.join(folder, f"preprocessed_{file_name}"), 'w', encoding='utf-8') as preprocessed_file:
40                     preprocessed_file.write(' '.join(preprocessed_text))
41

```

Graph Construction

Representation

We represented each document as a directed graph, where nodes represent unique terms (words) and edges denote term relationships based on their sequence in the text. The graph construction process involved the following steps:

- Reading the preprocessed text data from each document.
- Splitting the text into individual words to create nodes in the graph.
- Connecting adjacent words in the text with directed edges to represent the sequential order of terms.

Implementation

The graph construction was implemented in the `graph.py` module. We utilized the `networkx` library to create and manipulate graphs in Python. The `networkx.DiGraph` class was used to represent directed graphs, and edges were added between nodes to capture term relationships.

```

1 import os
2 import networkx as nx
3 import matplotlib.pyplot as plt
4
5 def construct_graph(file_path):
6     G = nx.DiGraph()
7     with open(file_path, 'r', encoding='utf-8') as file:
8         words = file.read().split()
9         for i in range(len(words) - 1):
10             current_word = words[i]
11             next_word = words[i + 1]
12             if G.has_edge(current_word, next_word):
13                 print(G[current_word][next_word])
14                 G[current_word][next_word]['weight'] += 1
15             else:
16                 G.add_edge(current_word, next_word, weight=1)
17     return G
18
19 def print_adjacency_matrix(G):
20     nodes = sorted(G.nodes())
21     print("\t" + "\t".join(nodes))
22     for node1 in nodes:
23         row = [str(G[node1].get(node2, {'weight': 0})['weight']) for node2 in nodes]
24         print(f"{node1}\t{'\t'.join(row)}")
25
26 def print_adjacency_list(G):
27     for node in G.nodes():
28         neighbors = list(G.neighbors(node))
29         if neighbors:
30             print(f"{node}: {' '.join(neighbors)}")
31         else:
32             print(f"{node}:")
33
34 def visualize_graph(G, file_path):
35     plt.figure(figsize=(10, 10))
36     pos = nx.spring_layout(G)
37     nx.draw(G, pos, with_labels=True, node_color='skyblue', node_size=10, edge_color='black', linewidths=1, arrows=True, arrowsize=20)
38     labels = nx.get_edge_attributes(G, 'weight')
39     nx.draw_networkx_edge_labels(G, pos, edge_labels=labels)
40     plt.title(f'Directed Graph of Words in {os.path.basename(file_path)}')
41     plt.show()
42
43 def visualize_preprocessed_files(folder):
44     graphs = []
45     for file_name in os.listdir(folder):
46         if file_name.startswith('preprocessed_') and file_name.endswith('.txt'):
47             file_path = os.path.join(folder, file_name)
48             G = construct_graph(file_path)
49             graphs.append(G)
50             # print_adjacency_matrix(G)
51             # print_adjacency_list(G)
52             # visualize_graph(G, file_path)
53     return graphs
54

```

Feature Extraction via Common Subgraphs

Identification

We employed frequent subgraph mining techniques to identify common subgraphs within the training set graphs. These common subgraphs served as features for classification, capturing the shared content across documents related to the same topic. The process involved:

- Generating candidate subgraphs from individual document graphs.
- Counting the frequency of each subgraph across all document graphs.
- Selecting the most frequent subgraphs as features for classification.

Implementation

The common subgraph identification was implemented in the `graph.py` module. We iterated through the training set graphs and used the `nx.compose` function from `networkx` to combine graphs and identify common subgraphs.

```

1  def find_common_subgraph(graphs):
2      if len(graphs) < 2:
3          print("At least two graphs are required to find a common subgraph.")
4          return None
5
6      common_subgraph = graphs[0].copy()
7      for graph in graphs[1:]:
8          common_subgraph = nx.compose(common_subgraph, graph)
9
10     # Remove nodes that are not common across all graphs
11     for node in list(common_subgraph.nodes()):
12         if not all(node in graph for graph in graphs):
13             common_subgraph.remove_node(node)
14
15     return common_subgraph
16
17 graphs=visualize_preprocessed_files(preprocess_files_folder)
18 # Find common subgraph
19 common_subgraph = find_common_subgraph(graphs)
20
21 if common_subgraph is not None:
22     print("Common subgraph found:")
23     print_adjacency_list(common_subgraph)
24     visualize_graph(common_subgraph, "Common Subgraph")
25 else:
26     print("No common subgraph found among the graphs.")
27
28 # Visualize preprocessed files as directed graphs
29

```

Classification with K-Nearest Neighbors (KNN)

This section describes the evaluation of the classification system using the K-Nearest Neighbors (KNN) algorithm. KNN is a non-parametric, lazy learning method that classifies data points based on the labels of their nearest neighbors in the feature space.

Code Implementation:

The provided code snippet demonstrates the implementation of KNN classification and evaluation metrics using the scikit-learn library. Key steps include:

1. Data Preprocessing:

- Splitting the data into training and testing sets using `train_test_split`.
- Converting the document data (`travelList`, `fashionList`, `diseaseList`) into feature vectors using the `graphsToVectors` function (assuming it transforms document graphs into numerical features).

2. KNN Model Creation:

- Initializing a `KNeighborsClassifier` object with `n_neighbors=3`, specifying the number of nearest neighbors to consider for classification.

3. Model Training:

- Training the KNN model on the training data (`x_train`, `y_train`).

4. Prediction on Test Set:

- Using the trained model to predict the class labels for the unseen test data (`x_test`).

5. Evaluation:

- Calculating accuracy, precision, recall, and F1-score using `accuracy_score` and `precision_recall_fscore_support` functions.
- Generating a confusion matrix using `confusion_matrix` to visualize the classification performance.
- The code also includes visualization of the confusion matrix using libraries like `matplotlib.pyplot`.

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.metrics import accuracy_score, precision_recall_fscore_support, confusion_matrix
4 from graph import returnTravel, returnFashion, returnDisease, graphToVector, graphsToVectors
5
6 travellist = returnTravel()
7 fashionList = returnFashion()
8 diseaselist = returnDisease()
9 docs = travellist + fashionList + diseaselist
10 labels = ['travel'] * len(travellist) + ['fashion'] * len(fashionList) + ['disease'] * len(diseaselist)
11
12 X_train, X_test, y_train, y_test = train_test_split(docs, labels, test_size=0.2, random_state=76)
13
14 # Initialize KNN classifier
15 X_train = graphsToVectors(X_train)
16 X_test = graphsToVectors(X_test)
17
18 classifier = KNeighborsClassifier(n_neighbors=3)
19
20 # Train the classifier
21 classifier.fit(X_train, y_train)
22
23 # Predict on the test set
24 predictions = classifier.predict(X_test)
25
```

Evaluation Metrics:

The following metrics were used to assess the performance of the KNN classifier:

- **Accuracy:** Measures the overall proportion of correctly classified data points.
- **Precision:** Represents the proportion of predicted positive labels that are actually positive.
- **Recall:** Captures the proportion of actual positive labels that were correctly predicted.
- **F1-score:** Combines precision and recall into a single metric, providing a balanced view of model performance.
- **Confusion Matrix:** Visualizes the distribution of actual vs. predicted labels, highlighting true positives, false positives, false negatives, and true negatives.

```

1 # Calculate accuracy
2 accuracy = accuracy_score(y_test, predictions)
3 print("Accuracy: ", int(accuracy*100))
4
5 # Calculate precision, recall, F1-score
6 precision, recall, f1_score, _ = precision_recall_fscore_support(y_test, predictions, average='weighted')
7
8 # Print precision, recall, and F1-score
9 print("Precision: ", int(precision*100))
10 print("Recall: ", int(recall*100))
11 print("F1-score: ", int(f1_score*100))
12
13 # Calculate the confusion matrix
14 cm = confusion_matrix(y_test, predictions)
15
16 # Import libraries for plotting (assuming you don't have them already)
17 import matplotlib.pyplot as plt
18
19 # Plot the confusion matrix
20 plt.figure(figsize=(8, 6))
21 plt.imshow(cm, interpolation='nearest', cmap=plt.cm.Blues)
22 plt.colorbar()
23
24 # Set labels for each axis
25 classes = ['travel', 'fashion', 'disease'] # Assuming these are your class labels
26 plt.xticks(range(len(classes)), classes, rotation=45)
27 plt.yticks(range(len(classes)), classes)
28
29 # Add text labels to each cell of the confusion matrix
30 thresh = cm.max() / 2.0
31 for i in range(len(cm)):
32     for j in range(len(cm)):
33         plt.text(j, i, cm[i, j], ha="center", va="center",
34                 color="white" if cm[i, j] > thresh else "black")
35
36 # Set labels for title and axes
37 plt.xlabel('Predicted Label')
38 plt.ylabel('True Label')
39 plt.title('Confusion Matrix')
40
41 # Show the confusion matrix plot
42 plt.tight_layout()
43 plt.show()
44

```

Results:

The provided code outputs the calculated values for accuracy, precision, recall, and F1-score. Additionally, the confusion matrix plot offers a detailed breakdown of the model's classification performance across different classes (travel, fashion, disease).

By analyzing these metrics, you can gain valuable insights into the effectiveness of the KNN classifier for your specific document classification task.

Results

Data Curation

The dataset was curated to ensure a balanced representation of three diverse topics: travel, fashion, and diseases. Each topic consisted of 15 pages of text, with approximately 300 words per page. This ensured that the dataset was adequately sized and representative of each class for effective model training and evaluation.

Clarity and Thoroughness of the Methodology

The methodology was meticulously defined and implemented, covering all essential steps of the document classification process:

- Data collection from various online sources
- Preprocessing of text data to prepare it for graph construction
- Graph representation of documents
- Feature extraction using common subgraphs

The clarity and thoroughness of the methodology facilitated a structured approach to model development and evaluation.

Creativity in Graph Representation and Feature Extraction

The use of graph-based features, particularly the identification of common subgraphs, showcased creativity in capturing the inherent relationships between terms within documents. By representing documents as directed graphs and extracting common subgraphs, the methodology leveraged graph theory principles to enhance the classification process. This creative approach allowed for a more nuanced understanding of document content beyond traditional vector-based models.

Depth of Analysis and Critical Reflection

Although not explicitly mentioned in the code, critical reflection on the challenges encountered during the project and the identification of potential improvements to the approach are essential aspects of any research endeavor. Through thoughtful analysis and reflection, insights into the strengths and limitations of the methodology can be gained, paving the way for future enhancements and contributions to the field of document classification.

Future Work

Integration of Advanced Graph-Based Techniques

In addition to KNN classification and common subgraph identification, future work will focus on integrating advanced graph-based techniques for document classification. This may include exploring graph neural networks (GNNs) and graph embedding methods to capture complex relationships and semantic information within document graphs. By leveraging these advanced techniques, the classification system can achieve greater accuracy and robustness across diverse document datasets.

Evaluation and Benchmarking

Future work will also involve comprehensive evaluation and benchmarking of the document classification system against existing state-of-the-art methods. This will include conducting comparative studies with traditional vector-based models and other graph-based approaches to assess the effectiveness and scalability of the proposed methodology. Through rigorous evaluation and benchmarking, the strengths and limitations of the classification system can be better understood, leading to informed decisions for further improvements.

Conclusion

In conclusion, this project has laid the foundation for a graph-based document classification system, encompassing data curation, graph construction, and feature extraction. While certain aspects such as KNN classification and common subgraph identification were not fully implemented due to a lack of understanding, the project has demonstrated creativity in graph representation and feature extraction. Moving forward, future work will focus on addressing the missing components and integrating advanced graph-based techniques to enhance the accuracy and efficiency of the classification system. Through continuous iteration and improvement, the proposed methodology holds promise for advancing the field of document classification and contributing to real-world applications in information retrieval and natural language processing.