

Kolapo Mogaji

Professor Mcmanus

ITAI 2376

February 26, 2025

Reflection Journal – Module 2, Lab 1: Processing Text

In this lab, I gained valuable insights into **natural language processing (NLP)** and its significance in handling text data. One of the most interesting aspects was creating a **word cloud**, which provided a visual representation of the most frequently occurring words in a dataset. This technique helped me understand how word frequency can reveal key themes within a body of text.

Another key takeaway was learning about **stemming and lemmatization**. I found it fascinating how stemming quickly reduces words to their root form, while lemmatization provides more linguistically accurate base words. Understanding the difference between the two is important because choosing the right approach depends on the specific needs of a project. For example, if speed is a priority, stemming might be preferred, whereas lemmatization would be better suited for applications requiring greater accuracy in word meaning.

Part-of-speech (POS) tagging was another powerful technique that enhanced my understanding of how words function within sentences. Seeing how words are categorized into nouns, verbs, adjectives, and other parts of speech reinforced the idea that language is structured and that computers can analyze this structure to extract meaning.

Additionally, exploring **Named Entity Recognition (NER)** demonstrated how machines can identify specific entities, such as people, places, and organizations. This is particularly useful in applications like information retrieval, customer sentiment analysis, and even fraud detection.

Overall, this lab reinforced the importance of **text preprocessing** as a critical step in NLP. Clean and well-processed text leads to better-performing models, whether in **sentiment analysis, chatbots, or recommendation systems**. This hands-on experience deepened my appreciation for the complexity of working with text data and highlighted how NLP techniques can transform raw text into meaningful insights.

Moving forward, I am eager to explore how these techniques can be integrated into **machine learning models** and applied to real-world scenarios, such as analyzing customer feedback or improving search engine accuracy.

Reflection Journal – Module 2, Lab 2: Using the BoW Method

In this lab, I explored the **Bag-of-Words (BoW) method**, which is a fundamental technique in natural language processing (NLP) for converting text into numerical representations. This method is crucial for preparing textual data for machine learning models. One of the most interesting aspects was understanding how BoW represents text by creating a vocabulary and encoding sentences as numerical vectors.

I started by working with **binary classification**, which simply records whether a word is present in a sentence without considering frequency. This method was useful for capturing word occurrences but lacked depth in understanding word importance. Moving forward, I learned about **word count-based BoW**, which includes frequency information, providing a more

informative representation. However, I realized that common words may dominate the data, making it harder to extract meaningful insights.

To address this, I explored **Term Frequency (TF)** and **Inverse Document Frequency (IDF)**.

TF normalizes word importance within a document, while IDF assigns lower weights to frequently used words and higher weights to rare words, making the representation more meaningful. Combining these two metrics using **TF-IDF** resulted in a more balanced approach to identifying significant words in text data.

One of the biggest takeaways from this lab was seeing how **vectorized text data can be used in machine learning applications**. I also realized that BoW methods, while simple, can result in **high-dimensional and sparse matrices** when working with large datasets, which could impact computational efficiency. This makes me curious about more advanced text representation techniques, such as **word embeddings** and **deep learning models for NLP**.

Overall, this lab deepened my understanding of how text data is processed for machine learning. The hands-on exercises reinforced key concepts, and I now feel more confident in applying BoW techniques to real-world NLP tasks, such as **sentiment analysis, document classification, and information retrieval**. Moving forward, I am eager to explore more sophisticated text vectorization methods to improve model performance.

Reflection Journal – GloVe Word Embeddings and Cosine Similarity

In this lab, I explored **GloVe word embeddings** and their application in understanding word relationships through **cosine similarity**. This experience deepened my understanding of how words can be represented as numerical vectors while maintaining their semantic meanings.

One of the most insightful parts of this lab was **loading pre-trained GloVe embeddings** and observing how words like **"computer"** and **"human"** have distinct yet meaningful vector representations. It was fascinating to see that these vectors are not just random numbers but structured in a way that captures linguistic and contextual relationships between words.

Another key takeaway was **using cosine similarity to measure word closeness**. By comparing words like **"car" and "truck" versus "car" and "bike,"** I could see how this mathematical approach helps quantify the semantic similarity between words. The results aligned with human intuition, reinforcing the effectiveness of GloVe embeddings in capturing real-world language relationships.

Additionally, I realized how **word embeddings play a crucial role in NLP applications**, including **chatbots, sentiment analysis, and recommendation systems**. Instead of working with raw text, machine learning models benefit from numerical representations of words, leading to better performance and interpretability.

A challenge I encountered was ensuring that the **correct pre-trained GloVe model was downloaded and loaded properly**. Managing large word vector files required careful file handling, and I learned the importance of verifying data paths and dimensions before running similarity calculations.

Overall, this lab provided a strong foundation in **word embeddings and similarity analysis**, preparing me for more advanced NLP tasks. Moving forward, I am excited to explore **other embedding techniques like Word2Vec and BERT** to compare their performance in various real-world scenarios.

Reflection on Improving the RNN Model

In this exercise, we built a simple **Recurrent Neural Network (RNN)** to process sequential data, specifically for text classification tasks. RNNs are powerful tools for sequential data because they can maintain an internal state (memory) across timesteps, allowing them to learn from the order of the data, which is crucial for tasks such as natural language processing, time-series analysis, and more.

While we have built a basic RNN, there are several ways we can **improve its performance** and make it more effective for complex tasks. Below are a few ways to enhance the model further:

1. Tuning Hyperparameters

- **Learning Rate:** The learning rate controls how fast or slow the model updates its weights during training. A **too high** learning rate can cause the model to overshoot the optimal solution, leading to poor convergence, while a **too low** learning rate can make the model converge too slowly. By experimenting with different learning rates (e.g., 0.001, 0.01, 0.0001), we can find the one that works best for our dataset.
- **Batch Size:** Batch size determines how many samples are processed at once during training. A **larger batch size** leads to more stable gradients but might require more memory, while a **smaller batch size** allows for faster updates but might have higher variance. Typically, values like **32**, **64**, and **128** are used for batch sizes.
- **Hidden Size:** The hidden size of the RNN determines the number of neurons in the hidden state. A **larger hidden size** means the model can capture more complex patterns,

but it will also require more computational resources and might lead to overfitting.

Testing different hidden sizes, such as **64**, **128**, or **256**, could help optimize the model.

2. Increasing the Number of Layers

RNNs can be enhanced by stacking multiple layers to make them **deeper**, which enables them to learn more abstract representations of the data. By increasing the `num_layers` parameter, the RNN will learn more complex features. For example, going from **1 layer** to **2 or 3 layers** might improve performance, especially for more complex tasks.

- **Num Layers:** You can experiment with `num_layers=2` or `num_layers=3` to stack multiple RNN layers. This helps the network capture richer temporal dependencies but may also make the training process more challenging.

3. Switching to GRU or LSTM

- **Gated Recurrent Units (GRU):** GRUs are an advanced type of RNN that aim to solve some of the problems inherent to vanilla RNNs, such as the **vanishing gradient problem**. GRUs use gates to control the flow of information and are computationally more efficient than LSTMs (though less powerful in some contexts). GRUs can be an excellent option if you need faster training with a decent level of performance.
 - [GRU Documentation](#)
- **Long Short-Term Memory (LSTM):** LSTMs are another advanced type of RNN, often considered superior for longer sequences and tasks that require learning from long-term dependencies. LSTMs are particularly useful when the model needs to remember

important information for long durations (e.g., in text, where the meaning of a word can depend on the context provided by previous words).

- [LSTM Documentation](#)

Both GRU and LSTM are widely used in NLP and time-series problems. Switching from a basic RNN to one of these models can potentially boost performance, especially when dealing with long sequences or complex data.

4. Further Experimentation

- **Regularization:** To prevent overfitting, especially as the model becomes more complex, you can incorporate **Dropout** between layers or use techniques like **Weight Regularization**.
- **Optimization Algorithms:** While SGD is commonly used, testing with more advanced optimizers like **Adam**, **RMSprop**, or **Adagrad** could result in faster convergence and better performance.
- **Sequence Padding and Masking:** For sequence data of varying lengths, you can experiment with **padding and masking** techniques to ensure the model only focuses on the real data without the padded values.

Conclusion

The flexibility of RNNs makes them a powerful tool for sequential data tasks. However, to improve model performance, you need to carefully experiment with different hyperparameters,

architectures, and types of RNNs. Increasing the number of layers, switching to GRU or LSTM, and fine-tuning parameters such as the learning rate and hidden size are effective strategies to enhance the model's ability to capture temporal dependencies and generalize better on unseen data.

By experimenting with these strategies, I can build a more robust model that performs better on complex tasks, particularly those involving long sequences of data such as text classification, time-series forecasting, and more.