

CAP 6640 – Natural Language Processing
Homework #4

Spring 2025
Due date: March 6, 2025

Problem 1: Describe how bidirectional RNNs work, including a diagram and mathematical formulation, and explain the problems they address.

Problem 2:

1. **Convolution** is the mathematical process of applying a kernel to an input to extract semantic or syntactic features from the given text. It involves element-wise multiplication followed by summation. In the context of NLP, convolutions are applied to word embeddings to detect n-gram patterns. For example, consider the sentence: "The concert was awesome.". First of all, we embed words into vectors:

$$\begin{aligned}\text{"The"} &\rightarrow [0.1, 0.2, 0.3] \\ \text{"concert"} &\rightarrow [0.2, 0.3, 0.4] \\ \text{"was"} &\rightarrow [0.3, 0.4, 0.5] \\ \text{"awesome"} &\rightarrow [0.4, 0.5, 0.6]\end{aligned}$$

Let's assume that we have a kernel of size 2, which is applied to the sentence. The kernel is represented as:

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.3 & 0.4 & 0.5 \end{bmatrix}$$

The convolution operation is performed as follows:

$$[0.1, 0.2, 0.3] \cdot [0.1, 0.2, 0.3] + [0.2, 0.3, 0.4] \cdot [0.3, 0.4, 0.5] = 0.1 \times 0.1 + 0.2 \times 0.2 + 0.3 \times 0.3 + 0.2 \times 0.3 + 0.3 \times 0.4 + 0.4 \times 0.5 = 0.01 + 0.04 + 0.09 + 0.06 + 0.12 + 0.2 = 0.52$$

2. **Padding** involves adding extra zeros or other relevant values around the input to control the output size and preserve sequence length in text. For example, without padding the sentence "The concert was awesome.", applying a kernel size of 3 reduces output length from 4 words to 2 words. With zero padding, we maintain the original length of the sentence as follows:

$$\begin{aligned}\phi &\rightarrow [0.0, 0.0, 0.0] \\ \text{"The"} &\rightarrow [0.1, 0.2, 0.3] \\ \text{"concert"} &\rightarrow [0.2, 0.3, 0.4] \\ \text{"was"} &\rightarrow [0.3, 0.4, 0.5] \\ \text{"awesome"} &\rightarrow [0.4, 0.5, 0.6] \\ \phi &\rightarrow [0.0, 0.0, 0.0]\end{aligned}$$

3. **Channels** are the number of feature maps generated by applying multiple kernels to the input. Each channel represents a different feature extracted from the text. This allows to combine different linguistic features. For example, if we apply 2 kernels to the sentence "The concert was awesome.", we get 2 channels, each representing a different feature. Assume that the other kernel is as follows:

$$\begin{bmatrix} 0.2 & 0.1 & 0.3 \\ 0.4 & 0.3 & 0.5 \end{bmatrix}$$

For the bigram "The concert", after applying the two kernels, we get two channels for representation as follows:

$$\text{"The", "concert"} \rightarrow [0.52, 0.50]$$

4. **Max Pooling** is a downsampling operation that selects the maximum value from a set of values in a given window while preserving key information. This operation helps to capture the most important features. For example, assume that we have the following feature map:

$$\begin{aligned}\phi, \text{"The"}, \text{"concert"} &\rightarrow [0.52, 0.50] \\ \text{"The"}, \text{"concert"}, \text{"was"} &\rightarrow [0.48, 0.56] \\ \text{"concert"}, \text{"was"}, \text{"awesome"} &\rightarrow [0.64, 0.42] \\ \text{"was"}, \text{"awesome"}, \phi &\rightarrow [0.40, 0.38]\end{aligned}$$

After applying max pooling for each channel separately, we get the following representation:

$$\max p \rightarrow [0.64, 0.56]$$

5. **Average Over Time** is a pooling operation that calculates the average of all feature activations across the sequence to capture the overall context. This operation helps to reduce the dimensionality of the feature map. Furthermore, averaging them might provide a sentence-level embedding. For example, assume that we have the following feature map:

$$\begin{aligned} \phi, \text{"The"}, \text{"concert"} &\rightarrow [0.52, 0.50] \\ \text{"The"}, \text{"concert"}, \text{"was"} &\rightarrow [0.48, 0.56] \\ \text{"concert"}, \text{"was"}, \text{"awesome"} &\rightarrow [0.64, 0.42] \\ \text{"was"}, \text{"awesome"}, \phi &\rightarrow [0.40, 0.38] \end{aligned}$$

After applying average over time for each channel separately, we get the following representation:

$$\text{avg} \rightarrow [0.51, 0.465]$$

6. **Striding** is the process of moving the kernel by a certain number of steps across the input during convolution to reduce the output size. This helps to control computational efficiency and feature map size. When we select stride as 1, the kernel moves one step at a time which outputs a dense feature map. For example, assume that we get the following feature map when we apply kernel size of 2 and stride of 1 to the sentence "The concert was awesome.":

$$\begin{aligned} \phi, \text{"The"}, \text{"concert"} &\rightarrow [0.52, 0.50] \\ \text{"The"}, \text{"concert"}, \text{"was"} &\rightarrow [0.48, 0.56] \\ \text{"concert"}, \text{"was"}, \text{"awesome"} &\rightarrow [0.64, 0.42] \\ \text{"was"}, \text{"awesome"}, \phi &\rightarrow [0.40, 0.38] \end{aligned}$$

For stride of 2 with the same kernel size, the feature map would be as follows:

$$\begin{aligned} \phi, \text{"The"}, \text{"concert"} &\rightarrow [0.52, 0.50] \\ \text{"concert"}, \text{"was"}, \text{"awesome"} &\rightarrow [0.64, 0.42] \end{aligned}$$

7. **K-Max Pooling** selects the top k largest activations values from the feature map instead of a fixed window. This retains significant features while preserving the order of the sequence. For example, assume that we have the following feature map:

$$\begin{aligned} \phi, \text{"The"}, \text{"concert"} &\rightarrow [0.52, 0.50] \\ \text{"The"}, \text{"concert"}, \text{"was"} &\rightarrow [0.48, 0.56] \\ \text{"concert"}, \text{"was"}, \text{"awesome"} &\rightarrow [0.64, 0.42] \\ \text{"was"}, \text{"awesome"}, \phi &\rightarrow [0.40, 0.38] \end{aligned}$$

After applying k-max pooling for each channel separately with k=2, we get the following representation:

$$\text{k-max p} \rightarrow \begin{bmatrix} 0.52 & 0.50 \\ 0.64 & 0.56 \end{bmatrix}$$

Problem 3:

A multi-layer RNN consist of multiple stacked RNN layers where the output of one layer is fed as input to the next layer. This hierarchical architecture allows the network to learn complex patterns and dependencies in the data across different levels of abstraction.

For an in input sequence $X = (x_1, x_2, \dots, x_T)$, a single-layer RNN updates its hidden state as:

$$h_t^{(1)} = f(W_h h_{t-1}^{(1)} + b_h)$$

where $h_t^{(1)}$ is the hidden state of the first layer at time step t , W_h is the weight matrix, b_h is the bias term, and f is the activation function.

What multi-layer RNNs do is to extend this formulation to multiple layers. In a multi-layer RNN, each hidden state at layer l takes input from the previous layer's hidden state as

$$h_t^{(l)} = f(W_h^{(l)} h_t^{(l-1)} + b_h^{(l)})$$

where $h_t^{(l)}$ is the hidden state of the l -th layer at time step t , $W_h^{(l)}$ is the weight matrix, $b_h^{(l)}$ is the bias term, and f is the activation function.

This enhances representation learning, allowing the network to learn hierarchical representations of the input data, where each layer processes information at a different level of abstraction, with lower layers capturing local patterns and deeper layers modeling high-level dependencies.

The advantages of using multiple layers in RNNs include:

1. **Hierarchical Representation Learning:** As mentioned earlier, multi-layer RNNs can learn hierarchical representations of the input data. Lower layers capture local patterns like words in a phrase, while deeper layers capture long-term dependencies like sentence meaning.
2. **Increased Model Capacity:** Compared to single layer RNNs, multi-layer RNNs can model complex sequential data better than shallow alternatives.
3. **Improved Generalization:** Hierarchical feature extraction of multi-layer RNNs makes the useful for complex sequences like a long text and speech.

However, multi-layer RNNs also face several challenges:

1. **Vanishing and Exploding Gradients:** Like in every RNN architecture, as gradients propagate through many layers and time steps, they become exponentially small, preventing earlier layers from learning.
2. **Long-Term Dependency:** Even though multi-layer RNNs are better than shallow ones, they still struggle to capture long-term dependencies in the data due to the vanishing gradient problem, which can limit their ability to model sequences with long-range dependencies.
3. **Computational Complexity and Cost:** Having multiple RNNs means more model parameters to train, which can require significant computational resources, making training slow.
4. **Optimization Challenges:** Training multi-layer RNNs can be challenging due to overfitting and convergence of gradients, which can make it harder to find an optimal model configuration.

Problem 4: Explain the encoder-decoder architecture, detailing its working mechanism and providing three real-world applications.

Problem 5: Explain how CNNs are applied to text understanding tasks, such as sentiment analysis.

Problem 6: Describe how neural machine translation operates and compare it to alternative approaches.

Problem 7:

Identify three limitations of RNNs that CNNs overcome.