

Issues with Stacking Self-attention layers

Name: Harin Raja Radha Krishnan

Suppose that we design a deep architecture to represent a sequence by stacking self-attention layers with positional encoding. What could be the issues?

When designing a deep architecture to represent sequences using stacked self-attention layers with positional encoding, several issues can arise. One concern is the **computational complexity** of the model. Stacking multiple self-attention layers can increase the computational requirements, making training and inference times longer, especially when dealing with large-scale datasets or lengthy sequences.

Another issue to consider is the potential for **overfitting**. Deep architectures with a large number of layers have a high capacity to memorize training examples, which can lead to poor generalization of unseen data. Regularization techniques such as dropout, weight decay, or early stopping may be necessary to mitigate overfitting and improve generalization performance.

The problem of **vanishing or exploding gradients** can also be encountered when stacking multiple self-attention layers. This can hinder the training process and make it difficult for the model to effectively learn from the data. Techniques like careful weight initialization, gradient clipping, or normalization methods such as batch normalization can be employed to stabilize the training process.

Interpretability is another important consideration when designing deep architectures. As the number of stacked layers increases, the model becomes less interpretable, making it challenging to understand the specific features or inputs that contribute to its predictions. This can be addressed by employing techniques like attention visualization or layer-wise relevance propagation to enhance the interpretability of the model.

Attention masking is another issue to address. Self-attention layers consider the entire input sequence when assigning attention weights, which means the model can attend to irrelevant or misleading information. Attention masking techniques can be used to

ensure that the model attends to the relevant aspects of the sequence while ignoring noise or irrelevant information, improving its performance and robustness.

Hyperparameter tuning is crucial for achieving optimal performance in deep architectures. Determining the number of self-attention layers, the dimensionality of positional encodings, or the learning rate can significantly impact the model's performance. Conducting thorough experiments and using techniques such as grid search or Bayesian optimization can help find the optimal configuration for the specific task at hand.

While these issues present challenges, they can be addressed through careful consideration and adaptation of deep learning techniques. By addressing computational complexity, overfitting, gradient instability, interpretability, attention masking, and hyperparameter tuning, it is possible to design more robust and effective deep architectures for representing sequences.