

MACHINE LEARNING

PART 35

Transformers in ML Explained

Before transformers, sequential models like RNNs and LSTMs were commonly used for natural language processing tasks. However, they had limitations in capturing long-range dependencies in sequences.

Transformers were introduced to address these issues and provide more efficient and effective sequence processing.

Key Components

- **Self-Attention Mechanism**

- Think of attention as a spotlight that the model can shine on different parts of a sequence.
- Self-attention allows the model to assign different levels of importance to different words in a sentence when making predictions.
- This mechanism is what enables transformers to process input data in parallel.

- **Multi-Head Attention**

- Instead of relying on a single attention mechanism, transformers use multiple heads, each attending to different parts of the input.
- This helps the model capture various aspects of the relationships in the data.

- **Positional Encoding**

- Since transformers don't inherently understand the order of the inputs, positional encoding is added to provide information about the positions of tokens in a sequence.

- **Feedforward Neural Networks**

- Each attention layer is followed by a feedforward neural network, which helps the model learn complex patterns and relationships.

Architecture

- **Encoder-Decoder Structure**
 - Transformers typically have an encoder-decoder structure.
 - The encoder processes the input sequence, and the decoder generates the output sequence.
- **Stacking Layers**
 - Both the encoder and decoder consist of multiple layers.
 - Each layer contains a self-attention mechanism and a feedforward neural network.

Training and Optimization

- **Loss Function**

- Commonly used loss functions include cross-entropy loss for classification tasks and mean squared error for regression tasks.

- **Backpropagation**

- The model is trained using backpropagation, adjusting the weights based on the error calculated by the loss function.

Applications

- **Natural Language Processing (NLP)**
 - Transformers excel in tasks like language translation, sentiment analysis, and text summarization.
- **Image Processing**
 - Transformers have been adapted for computer vision tasks, such as image classification and object detection.
- **Speech Recognition**
 - Transformers are used for processing sequential data in speech recognition systems.

Fine-Tuning and Transfer Learning

- *Pre-trained transformer models, such as BERT and GPT, are often fine-tuned on specific tasks for improved performance.*
- *Transfer learning with pre-trained models has become a common practice.*

Challenges and Future

- *Despite their success, transformers still face challenges, such as the need for large amounts of data and computational resources.*
- *Ongoing research aims to address these challenges and further enhance transformer models.*