

LN + Residual Connection

In the Transformer model, layer normalization is applied to the output of the self-attention layer before it is added to the residual connection. This helps stabilize and accelerate the training process by normalizing the inputs to each sub-layer.

Layer Normalization in Transformer

Layer normalization is applied in two key places within each Transformer encoder and decoder layer:

1. After Self-Attention:

- The output of the self-attention mechanism is passed through a layer normalization layer before being added to the residual connection (the original input to the self-attention layer).
- This helps in stabilizing the learning process by normalizing the inputs to the subsequent feed-forward network.

2. After Feed-Forward Network:

- Similarly, the output of the feed-forward network is also passed through a layer normalization layer before being added to its residual connection.

Mathematical Formulation

For the output of the self-attention layer, the layer normalization is applied as follows:

Let X be the input to the self-attention layer, and $\text{Attention}(X)$ be the output of the self-attention mechanism.

1. Residual Connection:

$$\text{Residual} = X + \text{LayerNorm}(\text{Attention}(X))$$

However, in the Transformer architecture, the layer normalization is actually applied after the self-attention and before the residual connection. The correct formulation is:

1. Layer Normalization:

$$\text{NormAttention} = \text{LayerNorm}(\text{Attention}(X))$$

2. Residual Connection:

$$\text{Output} = X + \text{NormAttention}$$

Layer Normalization Formula

Layer normalization computes the mean and variance across the features (i.e., the embedding dimension) of the input tensor. For an input tensor X with shape (N, T, D) , where N is the batch size, T is the sequence length, and D is the embedding dimension, layer normalization is applied as follows:

1. Compute Mean:

$$\mu = \frac{1}{D} \sum_{i=1}^D X_i$$

2. Compute Variance:

$$\sigma^2 = \frac{1}{D} \sum_{i=1}^D (X_i - \mu)^2$$

3. Normalize:

$$X_{\text{norm}} = \frac{X - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

where ϵ is a small value added for numerical stability.

4. Scale and Shift:

$$X_{\text{out}} = \gamma \cdot X_{\text{norm}} + \beta$$

where γ and β are learnable parameters.

Implementation of Layer Normalization

Here's a simplified implementation of layer normalization using NumPy:

```

1  import numpy as np
2
3  class LayerNormalization:
4      def __init__(self, embed_dim, epsilon=1e-6):
5          self.embed_dim = embed_dim
6          self.epsilon = epsilon
7          self.gamma = np.ones(embed_dim)
8          self.beta = np.zeros(embed_dim)
9
10     def forward(self, x):
11         # x shape: (N, T, D)
12         self.x = x
13         self.mean = np.mean(x, axis=-1, keepdims=True)
14         self.var = np.var(x, axis=-1, keepdims=True)
15         self.x_norm = (x - self.mean) / np.sqrt(self.var + self.epsilon)
16         return self.gamma * self.x_norm + self.beta
17
18     def backward(self, grad, lr):
19         # grad shape: (N, T, D)
20         N, T, D = grad.shape
21
22         # Gradient of loss with respect to gamma and beta
23         dL_dgamma = np.sum(grad * self.x_norm, axis=(0, 1))
24         dL_dbeta = np.sum(grad, axis=(0, 1))
25
26         # Update gamma and beta
27         self.gamma -= lr * dL_dgamma
28         self.beta -= lr * dL_dbeta
29
30         # Gradient of loss with respect to x_norm
31         dL_dx_norm = grad * self.gamma
32

```

```

33         # Gradient of loss with respect to x
34         dL_dx = dL_dx_norm / np.sqrt(self.var + self.epsilon)
35         dL_dx -= np.mean(dL_dx_norm, axis=-1, keepdims=True) / np.sqrt(self.var +
self.epsilon)
36         dL_dx -= (self.x - self.mean) * np.mean(dL_dx_norm * (self.x - self.mean), axis=-1,
keepdims=True) / (np.sqrt(self.var + self.epsilon) ** 3 * D)
37
38         return dL_dx

```

Usage Example

```

1  # Example usage:
2  embed_dim = 128
3  batch_size = 32
4  seq_len = 10
5
6  # Input tensor
7  x = np.random.randn(batch_size, seq_len, embed_dim)
8
9  # Initialize layer normalization
10 layer_norm = LayerNormalization(embed_dim)
11
12 # Forward pass
13 output = layer_norm.forward(x)
14 print("Output shape:", output.shape)
15
16 # Simulate loss gradient
17 grad = np.random.randn(batch_size, seq_len, embed_dim)
18
19 # Backward pass
20 dL_dx = layer_norm.backward(grad, lr=0.001)
21 print("Gradient shape:", dL_dx.shape)

```

Summary

Layer normalization in the Transformer model is applied to the output of the self-attention layer before it is added to the residual connection. This helps in stabilizing the training process by normalizing the inputs to each sub-layer. The layer normalization computes the mean and variance across the embedding dimension and normalizes the inputs using these statistics. The normalized inputs are then scaled and shifted using learnable parameters.

Residual connections

Residual connections, also known as skip connections, are a key component in deep neural networks, including Transformers. They help mitigate the vanishing gradient problem by allowing gradients to flow through the network more directly during backpropagation. Here's how they work and a simple implementation using NumPy.

How Residual Connections Work

1. Basic Concept:

- A residual connection wraps around a sub-layer (e.g., self-attention, feed-forward network).
- The input to the sub-layer is added to the output of the sub-layer.
- This allows the gradient to bypass the sub-layer directly, facilitating deeper network training.

2. Mathematical Formulation:

- Let X be the input to a sub-layer.
- Let $F(X)$ be the output of the sub-layer.
- The output of the residual connection is:

$$\text{Output} = X + F(X)$$

3. Layer Normalization:

- In Transformers, layer normalization is typically applied before the residual connection.
- The output of the sub-layer is normalized before being added to the input.

Implementation with NumPy

Below is a simple implementation of a residual connection with layer normalization:

```
1  import numpy as np
2
3  class ResidualConnection:
4      def __init__(self, embed_dim, epsilon=1e-6):
5          self.embed_dim = embed_dim
6          self.epsilon = epsilon
7          self.gamma = np.ones(embed_dim)
8          self.beta = np.zeros(embed_dim)
9
10     def forward(self, X, sublayer_output):
11         # X shape: (N, T, D)
12         # sublayer_output shape: (N, T, D)
13         self.X = X
14         self.sublayer_output = sublayer_output
15
16         # Compute mean and variance for layer normalization
17         self.mean = np.mean(sublayer_output, axis=-1, keepdims=True)
18         self.var = np.var(sublayer_output, axis=-1, keepdims=True)
19
20         # Normalize
```

```

21         self.sublayer_output_norm = (sublayer_output - self.mean) / np.sqrt(self.var +
self.epsilon)
22
23         # Scale and shift
24         self.sublayer_output_norm_scaled = self.gamma * self.sublayer_output_norm +
self.beta
25
26         # Residual connection
27         output = X + self.sublayer_output_norm_scaled
28         return output
29
30     def backward(self, grad, lr):
31         # grad shape: (N, T, D)
32         N, T, D = grad.shape
33
34         # Gradient of loss with respect to the output of the residual connection
35         dL_dOutput = grad
36
37         # Gradient of loss with respect to the scaled and shifted normalized sublayer output
38         dL_dSublayerOutputNormScaled = dL_dOutput
39
40         # Gradient of loss with respect to gamma and beta
41         dL_dGamma = np.sum(dL_dSublayerOutputNormScaled * self.sublayer_output_norm, axis=
(0, 1))
42         dL_dBeta = np.sum(dL_dSublayerOutputNormScaled, axis=(0, 1))
43
44         # Update gamma and beta
45         self.gamma -= lr * dL_dGamma
46         self.beta -= lr * dL_dBeta
47
48         # Gradient of loss with respect to the normalized sublayer output
49         dL_dSublayerOutputNorm = dL_dSublayerOutputNormScaled * self.gamma
50
51         # Gradient of loss with respect to the sublayer output
52         dL_dSublayerOutput = dL_dSublayerOutputNorm / np.sqrt(self.var + self.epsilon)
53         dL_dSublayerOutput -= np.mean(dL_dSublayerOutputNorm, axis=-1, keepdims=True) /
np.sqrt(self.var + self.epsilon)
54         dL_dSublayerOutput -= (self.sublayer_output - self.mean) *
np.mean(dL_dSublayerOutputNorm * (self.sublayer_output - self.mean), axis=-1, keepdims=True)
/ (np.sqrt(self.var + self.epsilon) ** 3 * D)
55
56         # Gradient of loss with respect to the input X
57         dL_dX = dL_dOutput
58
59         return dL_dX, dL_dSublayerOutput
60
61     # Example usage:
62     embed_dim = 128
63     batch_size = 32
64     seq_len = 10
65
66     # Input tensor
67     X = np.random.randn(batch_size, seq_len, embed_dim)
68
69     # Simulated output of a sub-layer (e.g., self-attention)
70     sublayer_output = np.random.randn(batch_size, seq_len, embed_dim)
71
72     # Initialize residual connection

```

```

73 residual_connection = ResidualConnection(embed_dim)
74
75 # Forward pass
76 output = residual_connection.forward(X, sublayer_output)
77 print("Output shape:", output.shape)
78
79 # Simulate loss gradient
80 grad = np.random.randn(batch_size, seq_len, embed_dim)
81
82 # Backward pass
83 dL_dX, dL_dSublayerOutput = residual_connection.backward(grad, lr=0.001)
84 print("Gradient shape (dL_dX):", dL_dX.shape)
85 print("Gradient shape (dL_dSublayerOutput):", dL_dSublayerOutput.shape)

```

Explanation

1. Forward Pass:

- The input X and the output of a sub-layer are passed through layer normalization.
- The normalized sub-layer output is scaled and shifted using learnable parameters γ and β .
- The residual connection adds the original input X to the scaled and shifted normalized sub-layer output.

2. Backward Pass:

- Gradients are computed for the loss with respect to the output of the residual connection.
- Gradients are backpropagated through the layer normalization and residual addition.
- Learnable parameters γ and β are updated using the computed gradients.

Summary

In the context of transformers and self-attention mechanisms, layer normalization (LayerNorm) is typically applied to the output of the self-attention layer (and also to the output of the feed-forward layer). The normalization is performed across the **feature dimension** (the embedding dimension) for each token independently.

Key Points:

1. Normalization Across Feature Dimension:

- For each token in the sequence, the mean and standard deviation are computed across the embedding features (the last dimension).
- This means that each token's features are normalized independently of other tokens in the sequence.

2. Formula:

- Given an input tensor $X \in \mathbb{R}^{n \times d}$, where n is the sequence length and d is the embedding dimension:

$$\text{LayerNorm}(X) = \frac{X - \text{mean}(X, \text{dim} = -1, \text{keepdim}=\text{True})}{\text{std}(X, \text{dim} = -1, \text{keepdim}=\text{True}) + \epsilon} \cdot \gamma + \beta$$

- Here:
 - $\text{mean}(X, \text{dim} = -1)$ computes the mean across the last dimension (features).
 - $\text{std}(X, \text{dim} = -1)$ computes the standard deviation across the last dimension.
 - γ and β are learnable parameters for scaling and shifting.
 - ϵ is a small value added for numerical stability.

3. Application in Transformers:

- Layer normalization is applied **after** the self-attention layer and **after** the feed-forward layer.
- It helps stabilize training by normalizing the distribution of inputs to each sub-layer.

Example Code for Layer Normalization

Here's a simple implementation of layer normalization:

```

1  import numpy as np
2
3  class LayerNorm:
4      def __init__(self, embedding_dim, epsilon=1e-5):
5          self.embedding_dim = embedding_dim
6          self.epsilon = epsilon
7          # Learnable parameters
8          self.gamma = np.ones(embedding_dim)
9          self.beta = np.zeros(embedding_dim)
10         # Gradients for learnable parameters
11         self.dgamma = np.zeros_like(self.gamma)
12         self.dbeta = np.zeros_like(self.beta)
13
14     def forward(self, x):
15         # Compute mean and std across the last dimension (features)
16         mean = np.mean(x, axis=-1, keepdims=True)
17         std = np.std(x, axis=-1, keepdims=True)
18         # Normalize
19         self.x_normalized = (x - mean) / (std + self.epsilon)
20         # Scale and shift
21         return self.gamma * self.x_normalized + self.beta
22
23     def backward(self, d_output):
24         # Gradient w.r.t. beta
25         self.dbeta = np.sum(d_output, axis=0)
26         # Gradient w.r.t. gamma
27         self.dgamma = np.sum(d_output * self.x_normalized, axis=0)
28         # Backpropagate gradient
29         d_x_normalized = d_output * self.gamma
30         d_x = d_x_normalized / (np.std(self.x_normalized, axis=-1, keepdims=True) +
self.epsilon)
31         return d_x

```

Explanation:

- **Forward Pass:**
 - Compute mean and standard deviation across the feature dimension.
 - Normalize the input using these statistics.
 - Scale and shift the normalized input using learnable parameters γ and β .
- **Backward Pass:**
 - Compute gradients for the learnable parameters γ and β .
 - Backpropagate the gradient through the normalization step.

Usage in Transformer Architecture:

Layer normalization is typically applied as follows in a transformer block:

1. After the self-attention layer:

$$\text{AttentionOutput} = \text{LayerNorm}(X + \text{SelfAttention}(X))$$

2. After the feed-forward layer:

$$\text{FinalOutput} = \text{LayerNorm}(\text{AttentionOutput} + \text{FeedForward}(\text{AttentionOutput}))$$

This normalization helps in stabilizing the training process and improving the convergence of the model.

Why Residual Connection?

Residual connections (or residual links) are a crucial component in transformer architectures, particularly in the context of self-attention and feed-forward layers. They serve several important purposes:

1. Addressing the Vanishing Gradient Problem

- **Issue:** In deep neural networks, gradients can vanish or explode during backpropagation, making it difficult to train deep models.
- **Solution:** Residual connections provide a direct path for gradients to flow through the network, bypassing one or more layers. This helps in maintaining gradient flow and mitigating the vanishing gradient problem.

2. Facilitating Training of Deep Networks

- **Issue:** Training very deep networks from scratch can be challenging due to degradation problems (where performance degrades as network depth increases).
- **Solution:** Residual connections make it easier to train deep networks by allowing the model to learn residual functions. Instead of learning to map inputs to outputs directly, the network learns to map inputs to the residual (difference) that needs to be added to the input. This is easier to optimize.

3. Stabilizing the Learning Process

- **Mechanism:** By adding the input of a layer to its output, residual connections provide a form of implicit identity mapping. This stabilizes the learning process and helps in maintaining the input's integrity through the network.

4. Enabling Skip Connections

- **Mechanism:** Residual connections allow the network to skip one or more layers. This is particularly useful in transformers where multiple self-attention and feed-forward layers are stacked.

5. Improving Information Flow

- **Mechanism:** Residual connections enhance the flow of information through the network. They allow features from earlier layers to be directly propagated to later layers, which can help in preserving important information across the network depth.

Residual Connections in Transformers

In a transformer architecture, residual connections are typically used around the self-attention layer and the feed-forward layer. Here's how they are applied:

1. After Self-Attention:

$$\text{AttentionOutput} = \text{LayerNorm}(X + \text{SelfAttention}(X))$$

2. After Feed-Forward Network:

$$\text{FinalOutput} = \text{LayerNorm}(\text{AttentionOutput} + \text{FeedForward}(\text{AttentionOutput}))$$

Why Are Residual Connections Necessary?

- **For Self-Attention:**

- Self-attention can sometimes produce outputs that are not well-aligned with the input. The residual connection ensures that the original input is retained and combined with the attention output, preventing the network from losing important information.

- **For Feed-Forward:**

- The feed-forward network processes each token independently. The residual connection helps in maintaining the context from the self-attention layer while allowing the feed-forward network to focus on local transformations.

Practical Benefits

- **Faster Convergence:** Models with residual connections generally converge faster during training.
- **Improved Performance:** They help in achieving better performance on tasks like machine translation, text generation, etc.
- **Simpler Optimization:** The training process becomes more stable and less sensitive to initialization.

In summary, residual connections are essential in transformers as they address several challenges associated with training deep networks, ensure stable gradient flow, and enhance the model's ability to learn complex patterns.

LN Further Explanation

In the original transformer architecture proposed in the paper "Attention Is All You Need" by Vaswani et al., layer normalization is applied in two places:

1. **After the self-attention layer.**
2. **After the feed-forward layer.**

This design choice is intentional and serves specific purposes. Let's break down why layer normalization is applied in both places and address the alternative viewpoints.

Why Layer Normalization is Applied in Both Places

1. **Stabilizing the Self-Attention Output:**

- The self-attention mechanism can produce outputs with varying scales and distributions, especially since it involves computing weighted sums of value vectors based on attention scores.
- Applying layer normalization after the self-attention layer helps stabilize the distribution of the attention outputs, making them more consistent and easier to process by subsequent layers.

2. **Stabilizing the Feed-Forward Output:**

- The feed-forward layer (typically a position-wise feed-forward network) applies non-linear transformations to each token's features independently.

- The outputs of this layer can also have varying scales and distributions. Layer normalization after the feed-forward layer helps maintain a consistent scale and distribution, which is beneficial for the following layers.

The Original Transformer Architecture

In the original transformer, each encoder and decoder layer consists of sub-layers (self-attention and feed-forward), and each sub-layer has a residual connection followed by layer normalization:

- **Encoder Layer:**

```
1 | Input -> Self-Attention -> Residual Connection -> LayerNorm -> Feed-Forward -> Residual  
   | Connection -> LayerNorm
```

- **Decoder Layer:**

```
1 | Input -> Self-Attention -> Residual Connection -> LayerNorm -> Cross-Attention -> Residual  
   | Connection -> LayerNorm -> Feed-Forward -> Residual Connection -> LayerNorm
```

Addressing the Alternative Viewpoints

Some people argue that layer normalization could be applied in only one of these places. However, the original transformer design and subsequent research suggest that applying it in both places is beneficial for the following reasons:

1. Improved Training Stability:

- Applying layer normalization after both the self-attention and feed-forward layers helps maintain stable gradients throughout the network. This is crucial for training deep models where gradient stability is a significant concern.

2. Consistent Feature Scales:

- Each sub-layer (self-attention and feed-forward) can introduce different scales and distributions to the data. Normalizing after each sub-layer ensures that the inputs to the next sub-layer are on a consistent scale, which helps the network learn more effectively.

3. Empirical Validation:

- The original transformer architecture has been extensively validated through experiments, and applying layer normalization in both places has proven to be effective in practice. Many subsequent transformer-based models (like BERT, GPT, etc.) have followed this design pattern.

The Truth

The original transformer architecture's approach of applying layer normalization after both the self-attention and feed-forward layers is well-founded and has been empirically validated. While it's theoretically possible to apply layer normalization in only one place, doing so would likely lead to less stable training and potentially suboptimal performance.

Practical Implementation

Here's a simplified representation of how layer normalization is applied in the transformer architecture:

```
1 class TransformerLayer:
2     def __init__(self, embedding_dim, num_heads, feed_forward_dim):
3         self.self_attention = SelfAttention(embedding_dim, num_heads)
4         self.layer_norm1 = LayerNorm(embedding_dim)
5         self.feed_forward = FeedForward(embedding_dim, feed_forward_dim)
6         self.layer_norm2 = LayerNorm(embedding_dim)
7
8     def forward(self, x):
9         # Self-attention sub-layer
10        attention_output = self.self_attention(x)
11        # Residual connection and layer normalization
12        x = self.layer_norm1(x + attention_output)
13
14        # Feed-forward sub-layer
15        ff_output = self.feed_forward(x)
16        # Residual connection and layer normalization
17        x = self.layer_norm2(x + ff_output)
18
19        return x
```

Key Takeaway

Applying layer normalization after both the self-attention and feed-forward layers is a core aspect of the transformer architecture. This design choice helps maintain stable gradients, consistent feature scales, and improved training dynamics. While alternative approaches might seem plausible, the original design has been extensively validated and is widely adopted in practice.

Feed-Forward Layer

The positional-wise feed-forward layer in the transformer model is a fully connected feed-forward network applied to each position in the sequence independently. It consists of two linear transformations with a ReLU activation function in between. Below is a detailed explanation and implementation:

Explanation of the Positional-Wise Feed-Forward Layer

The positional-wise feed-forward layer consists of two linear layers (dense layers) applied to each position in the sequence. The same linear layers are used for all positions, hence the name "positional-wise." The layer transforms the input at each position through a linear transformation, followed by a ReLU activation function, and then another linear transformation. This allows the model to learn non-linear relationships between the input features.

Implementation with NumPy

Here's a step-by-step implementation of the positional-wise feed-forward layer using NumPy:

```
1 import numpy as np
2
3 def initialize_weights(input_dim, hidden_dim, output_dim):
4     # Initialize weights for the two linear layers
5     W1 = np.random.randn(input_dim, hidden_dim) * np.sqrt(2.0 / input_dim)
6     b1 = np.zeros(hidden_dim)
7     W2 = np.random.randn(hidden_dim, output_dim) * np.sqrt(2.0 / hidden_dim)
8     b2 = np.zeros(output_dim)
9     return W1, b1, W2, b2
10
11 def positional_wise_feed_forward(x, W1, b1, W2, b2):
12     # First linear layer
13     hidden = np.dot(x, W1) + b1
14     # ReLU activation
15     hidden = np.maximum(hidden, 0)
16     # Second linear layer
17     output = np.dot(hidden, W2) + b2
18     return output
19
20 # Example usage
21 input_dim = 512
22 hidden_dim = 2048
23 output_dim = 512
24 seq_length = 10
25 batch_size = 32
26
27 # Generate random input data
28 x = np.random.randn(seq_length, batch_size, input_dim)
29
30 # Initialize weights
31 W1, b1, W2, b2 = initialize_weights(input_dim, hidden_dim, output_dim)
32
33 # Forward pass
34 output = positional_wise_feed_forward(x, W1, b1, W2, b2)
35
36 print(output.shape)
```

Key Points

- Positional Independence:** The same weights are used for all positions in the sequence, allowing the layer to treat each position independently while sharing parameters across positions.
- Non-Linear Transformation:** The ReLU activation function introduces non-linearity, enabling the model to capture complex relationships between input features.
- Dimensionality Expansion and Reduction:** The hidden layer typically expands the dimensionality of the input to a larger value (e.g., 2048) and then reduces it back to the original dimension (e.g., 512), allowing the model to learn richer representations.

The positional-wise feed-forward layer plays a crucial role in the transformer architecture by adding non-linear transformations to the output of the multi-head attention layer, enhancing the model's ability to capture complex patterns in the data.