

Backpropagation in Convolutional Neural Networks: Engineering Challenges

Lewis Levin and ChatGPT

March 22, 2005

Backpropagation in Convolutional Neural Networks: Engineering Challenges

Backpropagation in convolutional neural networks (CNNs) involves several intricate engineering challenges that are often glossed over in theoretical explanations. This report details the specific computational complexities of backpropagation between convolutional layers, using a typical MNIST digit recognition network as our example.

The Challenge of Channel Dimensionality Transformation

When backpropagating from layer 3 to layer 2 in our example network, we face the non-trivial task of transforming the gradient signal from $24 \times 24 \times 8 \times N$ to $26 \times 26 \times 4 \times N$. This transformation involves not only expanding the spatial dimensions but also reducing the number of channels, which creates unique implementation challenges.

Spatial Dimension Expansion

The spatial dimension expansion (from 24×24 to 26×26) follows established principles in convolutional backpropagation. According to the convolution backward process, we need to pad the incoming gradient with zeros before applying the convolution operation¹. This padding is typically $(\text{kernel_size} - 1)$ on each side of the feature map, which accounts for the reduction in spatial dimensions that occurred during the forward pass².

For the backpropagation process, we must:

1. Pad the gradient coming from layer 3 (dL/dZ) with zeros
2. Flip the filters of layer 3 (as required by the mathematics of convolution)³⁴
3. Perform the backward convolution operation with these flipped filters

The padding ensures that the output of the backward convolution has the correct spatial dimensions that match layer 2's activation map⁵.

Channel Reduction: The Core Challenge

The more complex aspect of this backpropagation is reducing the number of channels from 8 back to 4. This is challenging because the forward pass expanded from 4 to 8 channels through the use of 8 different filters, each with 4 input channels⁶.

The key insight is that during backpropagation, the gradient for each input channel in layer 2 is computed by summing contributions from all output channels in layer 3 that were affected by that input channel⁷. This summation naturally handles the reduction from 8 channels back to 4.

¹https://deeplearning.cs.cmu.edu/F22/document/homework/HW2/HW2P1_F22.pdf

²<https://deeplearning.cs.cmu.edu/F22/document/slides/lec12.CNN4.pdf>

³https://deeplearning.cs.cmu.edu/F22/document/homework/HW2/HW2P1_F22.pdf

⁴<http://liushuaicheng.org/TIP/SlimConv.pdf>

⁵https://deeplearning.cs.cmu.edu/F22/document/homework/HW2/HW2P1_F22.pdf

⁶https://deeplearning.cs.cmu.edu/F22/document/homework/HW2/HW2P1_F22.pdf

⁷<https://www.jefkine.com/general/2016/09/05/backpropagation-in-convolutional-neural-networks/>

Mathematically, for each input channel m in layer 2, the gradient is calculated as:

$$dL/dA[:, :, m, :] = \text{sum_over_k}(\text{convolution}(dL/dZ[:, :, k, :], \text{flipped_W}[:, :, m, k]))$$

Where:

- dL/dZ is the gradient with respect to layer 3's pre-activation output
- W is the weight matrix of layer 3's convolution
- k iterates over all 8 output channels of layer 3⁸⁹

Applying the Leaky ReLU Derivative

Before backpropagating through the convolution operation itself, we must apply the derivative of the leaky ReLU activation function to the incoming gradient. The derivative of leaky ReLU is straightforward:

- 1 for inputs > 0
- A small positive value (typically 0.01) for inputs ≤ 0 ¹⁰

This derivative is applied element-wise to the gradient signal before performing the convolution operation:

$$dL/dZ = dL/dA * f'(Z)$$

Where $f'(Z)$ is the derivative of leaky ReLU evaluated at the pre-activation values of layer 2¹¹.

Implementation Approach

To implement this backpropagation process correctly, we would:

1. Start with the gradient dL/dZ (dimensions: $24 \times 24 \times 8 \times N$)
2. Apply the leaky ReLU derivative based on the cached pre-activation values Z
3. Initialize dL/dA with zeros (dimensions: $26 \times 26 \times 4 \times N$)
4. For each output channel k in layer 3 ($k=1,2,\dots,8$):
 - For each input channel m in layer 2 ($m=1,2,3,4$):
 - Pad $dL/dZ[:, :, k, :]$ with zeros to match required dimensions
 - Flip $W[:, :, m, k]$ (the filter connecting channel m in layer 2 to channel k in layer 3)
 - Perform convolution between padded $dL/dZ[:, :, k, :]$ and flipped $W[:, :, m, k]$
 - Add the result to $dL/dA[:, :, m, :]$ ¹²¹³
5. Apply the derivative of leaky ReLU to get dL/dZ

This procedure ensures that the gradient signal is correctly transformed from $24 \times 24 \times 8 \times N$ to $26 \times 26 \times 4 \times N$, preserving the chain rule of calculus throughout the network.

Practical Considerations

Several practical considerations further complicate this implementation:

Memory Management

The backpropagation process requires caching the pre-activation values (Z values) for each layer during the forward pass¹⁴. For large networks, this can consume significant memory.

⁸<https://www.jefkine.com/general/2016/09/05/backpropagation-in-convolutional-neural-networks/>

⁹https://www.reddit.com/r/deeplearning/comments/ndmoau/convolutional_layer_multichannel_backpropagation/

¹⁰<https://stackoverflow.com/questions/32546020/neural-network-backpropagation-with-relu>

¹¹<https://stackoverflow.com/questions/32546020/neural-network-backpropagation-with-relu>

¹²https://deeplearning.cs.cmu.edu/F22/document/homework/HW2/HW2P1_F22.pdf

¹³<https://www.jefkine.com/general/2016/09/05/backpropagation-in-convolutional-neural-networks/>

¹⁴<https://stackoverflow.com/questions/32546020/neural-network-backpropagation-with-relu>

Computational Efficiency

The naive implementation described above involves multiple convolution operations ($8 \times 4 = 32$ convolutions in our example). In practice, these operations can be parallelized and optimized using techniques such as grouped convolutions¹⁵ or implementing the operations as matrix multiplications.

Framework Implementation

Most deep learning frameworks abstract away these details, but understanding them is crucial for debugging, optimization, and developing custom layers. Frameworks like TensorFlow and PyTorch internally implement these operations efficiently using highly optimized linear algebra libraries¹⁶.

Conclusion

Backpropagation in convolutional neural networks, particularly the channel reduction aspect, represents a significant engineering challenge that is often overlooked in theoretical discussions. The process involves careful management of dimensions, proper application of the chain rule, and efficient implementation of multiple convolution operations.

By understanding these challenges in detail, practitioners can better debug their networks, optimize performance, and develop custom architectures with confidence. The example presented here, focusing on backpropagation from a $24 \times 24 \times 8 \times N$ layer to a $26 \times 26 \times 4 \times N$ layer, illustrates the intricate nature of gradient flow in convolutional architectures and highlights the importance of proper implementation details in deep learning systems.

¹⁵<https://stackoverflow.com/questions/68841748/backprop-for-repeated-convolution-using-grouped-convolution>

¹⁶<https://www.youtube.com/watch?v=Betg6UI9d0Q>