



INFO 7375, Neural Networks & AI

Home Work to Chapter - 15

Submitted By:

Abdul Haseeb Khan

NUID: 002844724

khan.abdulah@northeastern.edu

What is spatial separable convolution and how is it different from simple convolution?

Spatial Separable Convolution

Spatial separable convolution is a technique that decomposes a 2D convolution operation into two sequential 1D convolutions, one applied along the horizontal dimension and another along the vertical dimension. This approach fundamentally changes how the convolution is computed compared to standard convolution, offering computational advantages at the cost of reduced expressiveness.

In standard convolution, a 2D kernel (such as a 3×3 filter) is applied directly to the input image or feature map. Each output value is computed by taking the element-wise multiplication of the kernel weights with the corresponding input values and summing them all together. For a 3×3 kernel, this means 9 multiplications and additions per output pixel, using 9 independent parameters that can represent any pattern within that 3×3 window.

Spatial separable convolution takes a different approach by factorizing the 2D kernel into two 1D kernels. First, a horizontal kernel (like a 1×3 filter) convolves across each row of the input, producing an intermediate result. Then, a vertical kernel (like a 3×1 filter) convolves down each column of that intermediate result to produce the final output. Mathematically, this is equivalent to using a 2D kernel that is the outer product of the two 1D kernels. For a 3×3 effective receptive field, you only need 6 parameters (3 for horizontal, 3 for vertical) instead of 9, and only 6 multiplications per output pixel instead of 9.

The computational efficiency becomes even more pronounced with larger kernels. For a $k \times k$ kernel, standard convolution requires k^2 operations and parameters, while spatial separable convolution only needs $2k$ operations and parameters. This means for a 5×5 kernel, you reduce from 25 operations to just 10, a 60% reduction in computation. This efficiency made spatial separable convolution particularly valuable in traditional image processing applications where large kernels were common for operations like Gaussian blurring.

However, this efficiency comes with a significant limitation in representational power. Not all 2D kernels can be decomposed into separable form, only those that have rank 1, meaning they can be expressed as the outer product of two vectors. While kernels for operations like Gaussian blur, box blur, and simple horizontal or vertical edge detection (like Sobel filters) are naturally separable, many important kernels cannot be separated. Arbitrary rotation kernels, diagonal edge detectors at non-standard angles, and complex texture detectors all require the full expressiveness of a 2D kernel. This is why most learned kernels in modern CNNs are not separable, the training process typically discovers patterns that require the full flexibility of standard convolution.

What is the difference between depth-wise and pointwise convolutions?

Depthwise Convolution applies a single convolutional filter per input channel. For an input with C channels, it uses C different $k \times k$ kernels, where each kernel convolves with only its corresponding channel. The operation preserves the number of channels (C input channels $\rightarrow C$ output channels) and performs spatial filtering without mixing channel information. The parameter count is $C \times k \times k$.

Pointwise Convolution is a 1×1 convolution that operates across channels at each spatial location. It performs a linear combination of all input channels to produce each output channel, effectively mixing channel information without any spatial filtering. For C input channels and M output channels, it uses $C \times M$ parameters.

Key Differences:

The primary distinction is the dimension of operation. Depthwise convolution performs spatial filtering independently on each channel (operates in the spatial domain), while pointwise convolution performs channel mixing at each spatial position (operates in the channel domain). Depthwise convolution cannot combine information across channels, and pointwise convolution cannot capture spatial patterns.

Computational Complexity:

- Standard $k \times k$ convolution: $C \times M \times k \times k$ parameters
- Depthwise $k \times k$ convolution: $C \times k \times k$ parameters
- Pointwise convolution: $C \times M$ parameters

Depthwise Separable Convolution combines both operations sequentially: depthwise convolution first (spatial filtering per channel), followed by pointwise convolution (channel mixing). This factorization reduces parameters from $C \times M \times k \times k$ to $(C \times k \times k) + (C \times M)$, achieving significant computational savings while maintaining reasonable representational capacity.

For example, a 3×3 standard convolution with 32 input and 64 output channels requires 18,432 parameters. The depthwise separable version requires only $288 + 2,048 = 2,336$ parameters, an $8 \times$ reduction. This efficiency-expressiveness trade-off makes depthwise separable convolution the foundation of efficient CNN architectures like MobileNet and EfficientNet.

What is the sense of 1×1 convolution?

What it does: A 1×1 convolution applies a linear transformation across channels at each spatial position independently. Despite having no spatial extent, it performs a weighted sum of all input channels to produce each output channel value at every pixel location.

1. Channel Dimension Reduction/Expansion

1×1 convolutions efficiently change the number of channels without affecting spatial dimensions. This is crucial for controlling computational cost - reducing channels before expensive operations (like 3×3 or 5×5 convolutions) and expanding them afterward. This forms the basis of bottleneck architectures in ResNet and Inception modules.

2. Cross-Channel Information Mixing

While standard spatial convolutions (3×3 , 5×5) primarily capture spatial patterns, 1×1 convolutions learn relationships between channels. They can learn cross-channel

correlations and feature combinations, essentially performing feature transformation in the channel dimension.

3. Adding Nonlinearity

When followed by activation functions (ReLU, etc.), 1×1 convolutions add nonlinear transformations without changing spatial structure. This increases network depth and representational power with minimal computational overhead.

4. Linear Projection

1×1 convolutions act as learnable linear projections, similar to fully connected layers but preserving spatial structure. Each spatial position gets the same transformation, making them position-independent feature projectors.

Computational Efficiency:

For C input channels and M output channels, a 1×1 convolution requires only $C \times M$ parameters and $C \times M$ multiplications per spatial position. This is significantly cheaper than spatial convolutions while still providing learnable transformations.

Common Use Cases:

- **Network-in-Network (NiN):** Introduced 1×1 convolutions for deeper local modeling
- **Inception/GoogLeNet:** Dimension reduction before expensive convolutions
- **ResNet:** Matching dimensions in skip connections
- **Depthwise Separable Convolution:** Channel mixing after depthwise spatial filtering
- **Feature Pyramid Networks:** Channel alignment across different scales

The key insight is that 1×1 convolution provides a computationally efficient way to transform the channel dimension while maintaining spatial structure, making it an essential building block in modern CNN architectures.

What is the role of residual connections in neural networks?

Residual connections (skip connections) create shortcuts that allow the input of a layer to bypass one or more layers and be added directly to the output of a deeper layer. Mathematically, instead of learning $F(x)$, the network learns $F(x) + x$, where x is the input and $F(x)$ is the residual function.

Primary Benefits:

1. Solving the Vanishing Gradient Problem

In deep networks, gradients often diminish exponentially during backpropagation, making it difficult to train early layers. Residual connections provide a direct gradient pathway from deeper layers to shallower ones. The gradient flow includes an identity term (derivative of x is 1), ensuring at least some gradient reaches early layers even if $F(x)$ produces small gradients.

2. Enabling Very Deep Networks

Before ResNet, networks deeper than 20-30 layers often performed worse than shallower ones - not due to overfitting but optimization difficulty. Residual connections allow training of networks with hundreds or even thousands of layers by making the optimization landscape smoother and easier to navigate.

3. Identity Mapping as Default

The network only needs to learn the residual $F(x) = 0$ to implement an identity mapping (output = input). This is easier than learning the full identity function from scratch. If additional layers aren't needed, the network can effectively "skip" them by driving residual weights toward zero.

4. Feature Reuse and Gradient Highways

Residual connections enable features from earlier layers to be directly accessible to deeper layers, promoting feature reuse throughout the network. They create multiple paths of different lengths through the network, similar to an ensemble of shallow and deep networks.

Mathematical Formulation:

- Standard layer: $y = F(x, W)$
- Residual block: $y = F(x, W) + x$
- With dimension change: $y = F(x, W) + Wx$ (where W is a linear projection)

Gradient Flow Analysis:

During backpropagation, the gradient includes: $\partial \text{loss} / \partial x = \partial \text{loss} / \partial y \times (\partial F / \partial x + 1)$. The "+1" term ensures gradient flow even when $\partial F / \partial x$ approaches zero, preventing gradient vanishing.

Impact on Training Dynamics:

- **Faster convergence:** Networks with residual connections typically converge faster than their plain counterparts
- **Better initialization:** The network starts closer to identity mapping, a reasonable initial state
- **Smoother loss landscape:** Residual connections create smoother optimization surfaces with fewer poor local minima

Common Applications:

- **ResNet:** Original application, enabling 152+ layer networks
- **DenseNet:** Extends the concept with dense connections to all subsequent layers
- **Transformer architectures:** Essential in attention blocks and feed-forward networks
- **U-Net:** Skip connections between encoder and decoder for preserving spatial information

The fundamental insight is that learning residual functions (the difference from identity) is easier than learning arbitrary transformations from scratch, especially as networks become deeper. This simple modification revolutionized deep learning by making truly deep networks trainable and practical.