

CHAPTER 16

Deep learning

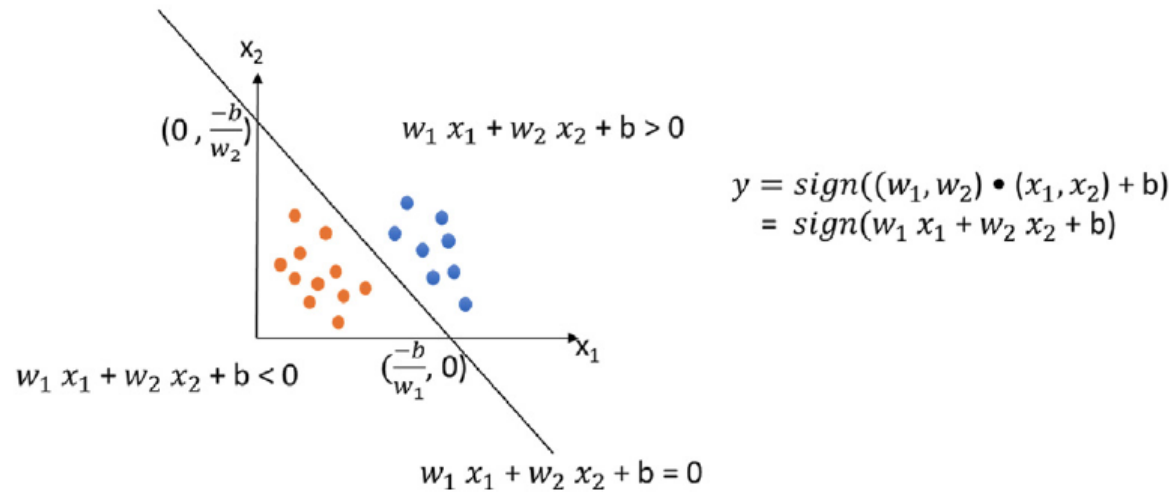


FIGURE 16.1

A perceptron linear classifier example in which the input is a two-dimensional vector.

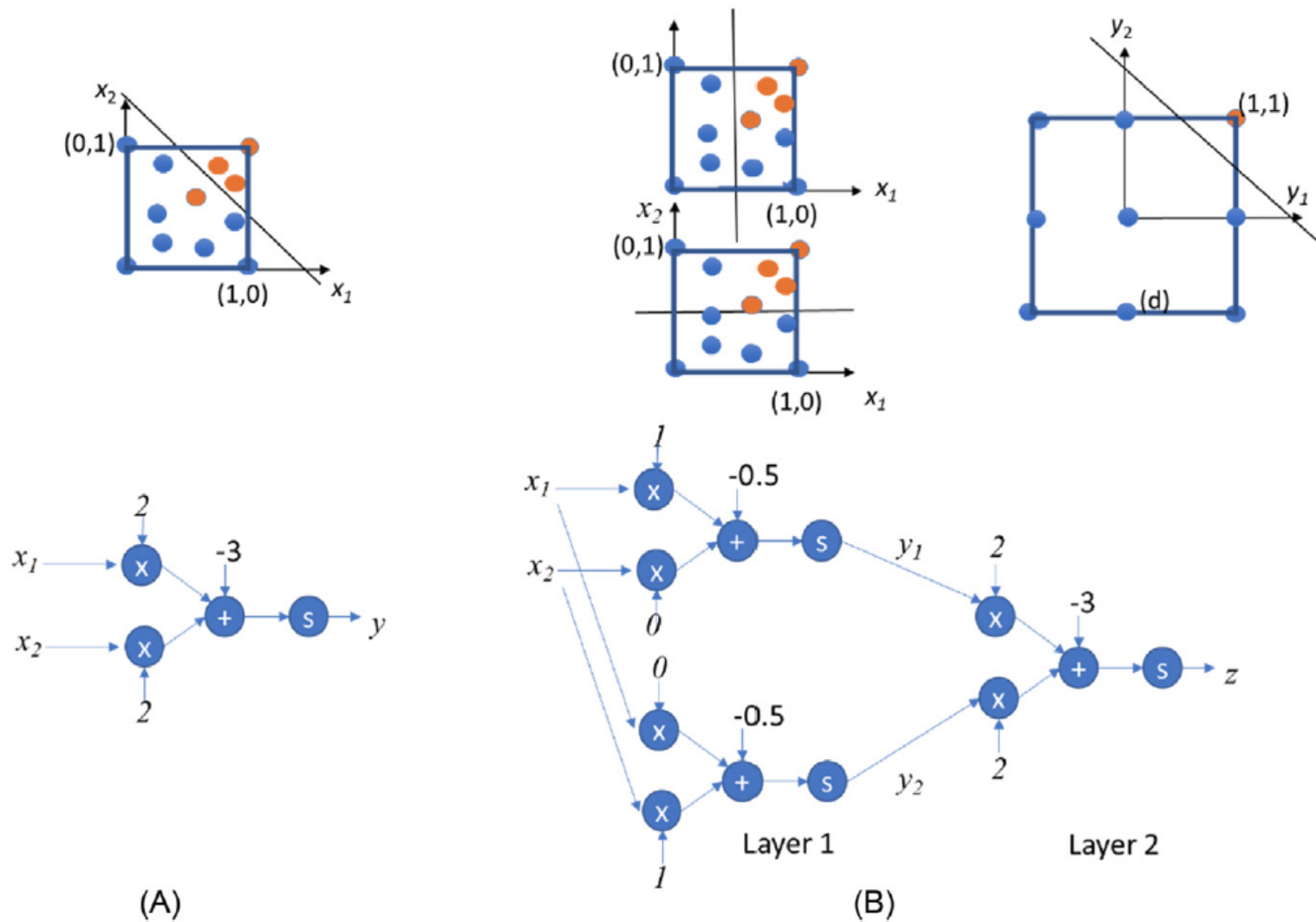


FIGURE 16.2

A multilayer perceptron example.

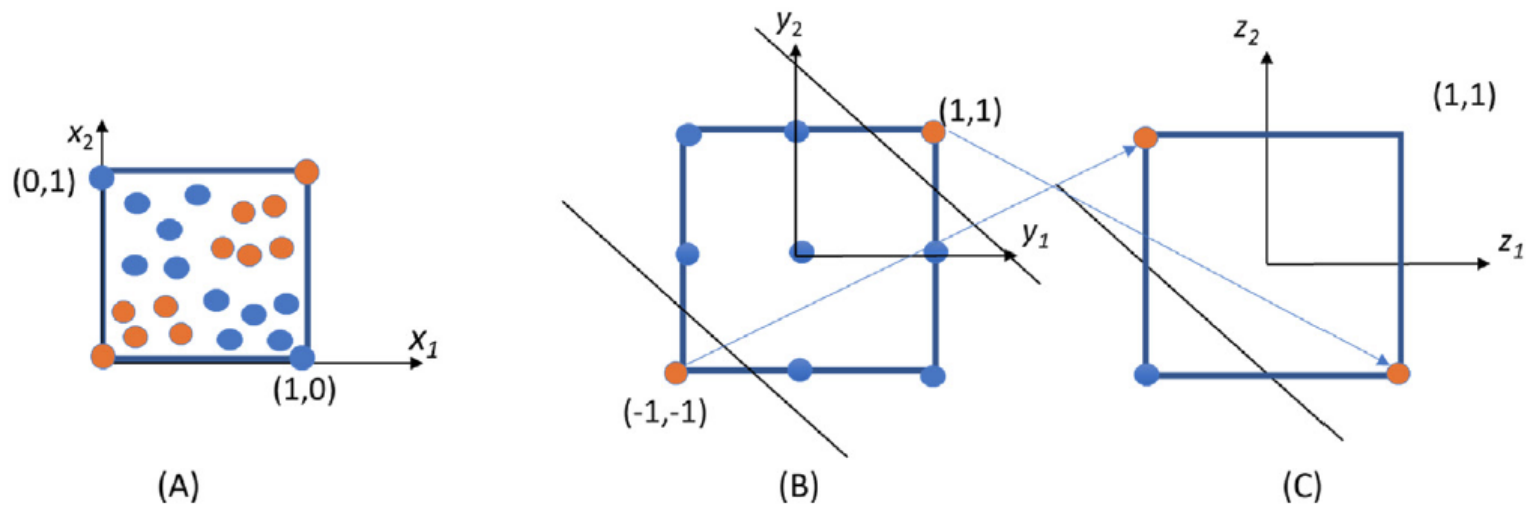


FIGURE 16.3

Need for perceptrons with more than two layers.

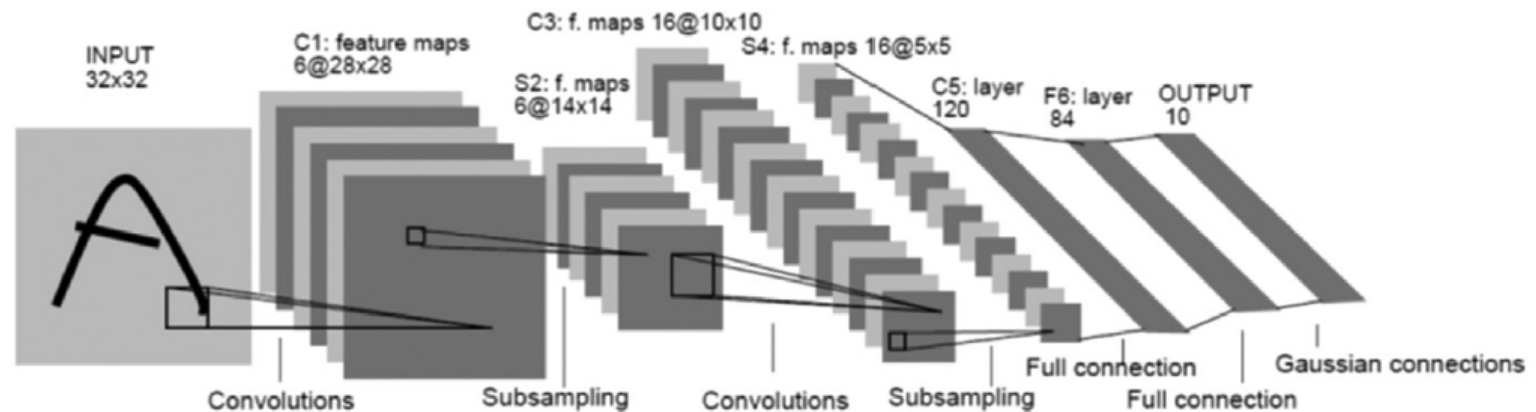


FIGURE 16.4

LeNet-5, a convolutional neural network for handwritten digit recognition. The letter A in the input should be classified as none of the ten classes (digits).

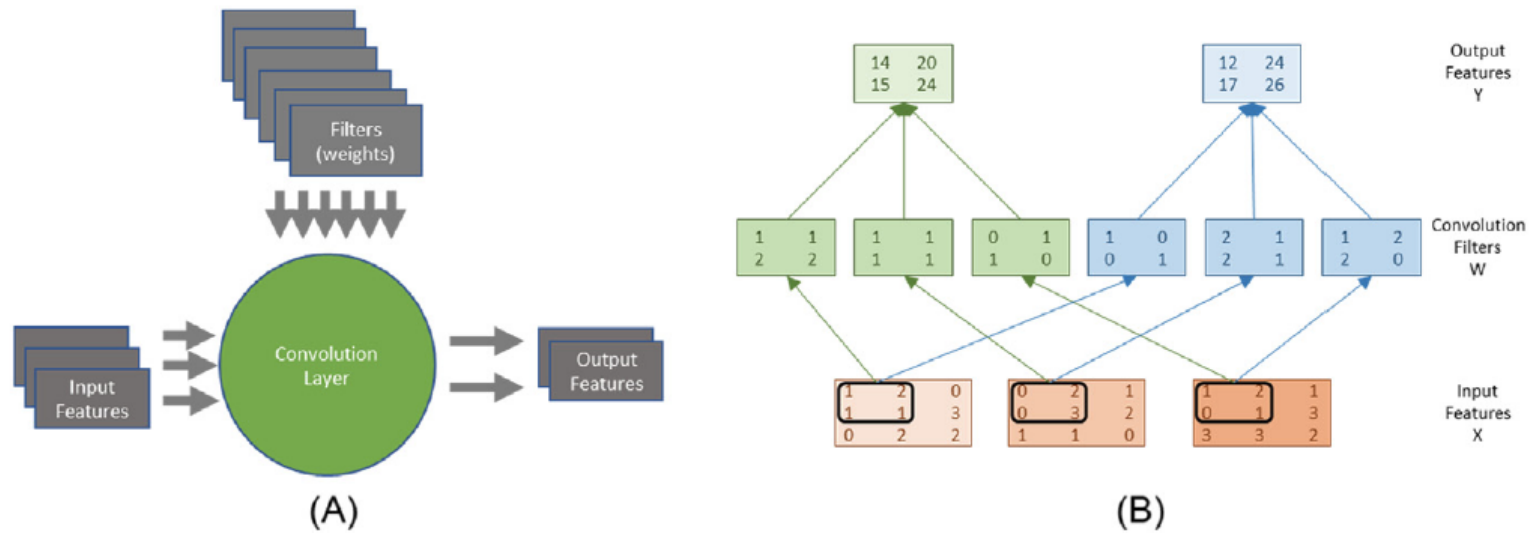


FIGURE 16.5

Forward propagation path of a convolutional layer.

```

01 void convLayer_forward(int M, int C, int H, int W, int K, float* X, float* W,
    float* Y) {
02     int H_out = H - K + 1;
03     int W_out = W - K + 1;
04     for(int m = 0; m < M; m++)          // for each output feature map
05         for(int h = 0; h < H_out; h++)    // for each output element
06             for(int w = 0; w < W_out; w++) {
07                 Y[m, h, w] = 0;
08                 for(int c = 0; c < C; c++) // sum over all input feature maps
09                     for(int p = 0; p < K; p++) // KxK filter
10                         for(int q = 0; q < K; q++)
11                             Y[m, h, w] += X[c, h + p, w + q] * W[m, c, p, q];
12             }
13 }

```

FIGURE 16.6

A C implementation of the forward propagation path of a convolutional layer.

```

01 void subsamplingLayer_forward(int M, int H, int W, int K, float* Y, float*
S){
02     for(int m = 0; m < M; m++)                // for each output feature map
03         for(int h = 0; h < H/K; h++)            // for each output element,
04             for(int w = 0; w < W/K; w++) {      // this code assumes that H and W
05                 S[m, x, y] = 0.;                // are multiples of K
06                 for(int p = 0; p < K; p++) {    // loop over KxK input samples
07                     for(int q = 0; q < K; q++)
08                         S[m, h, w] += Y[m, K*h + p, K*w+ q] / (K*K);
09                 }
10                 // add bias and apply non-linear activation
11                 S[m, h, w] = sigmoid(S[m, h, w] + b[m]);
12     }

```

FIGURE 16.7

A sequential C implementation of the forward propagation path of a subsampling layer. The layer also includes an activation function, which is included in a convolutional layer if there is no subsampling layer after the convolutional layer.

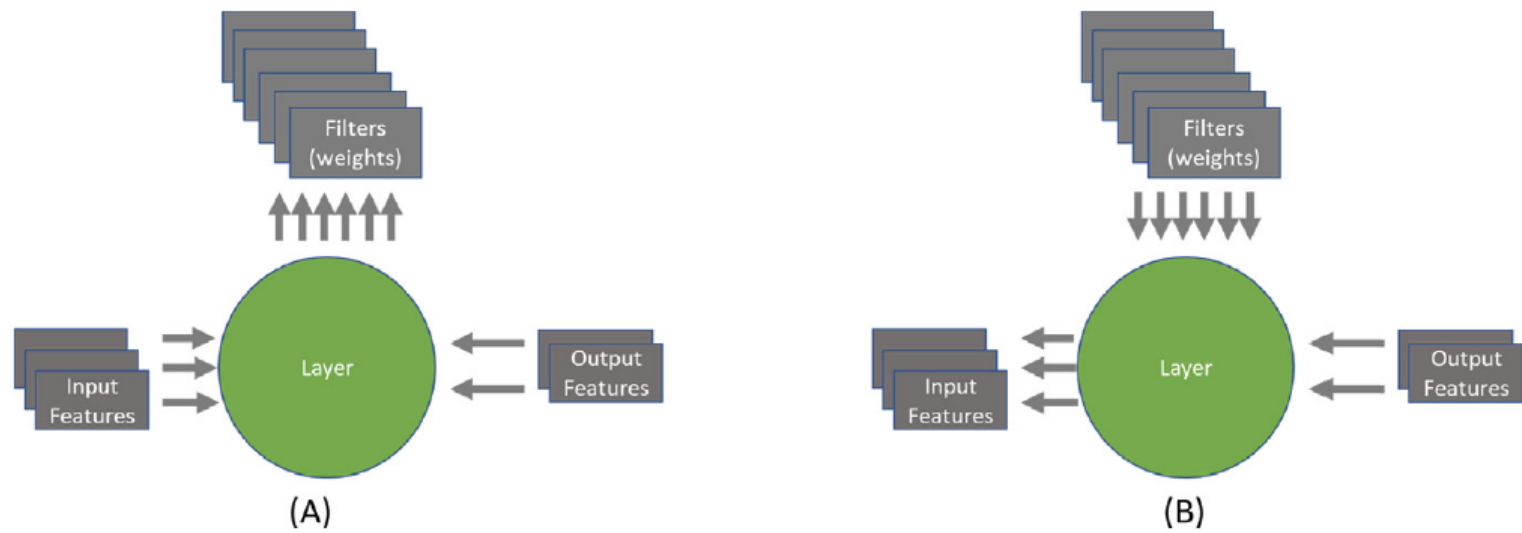


FIGURE 16.8

Backpropagation of (A) $\frac{\partial E}{\partial w}$ and (B) $\frac{\partial E}{\partial x}$ for a layer in CNN.

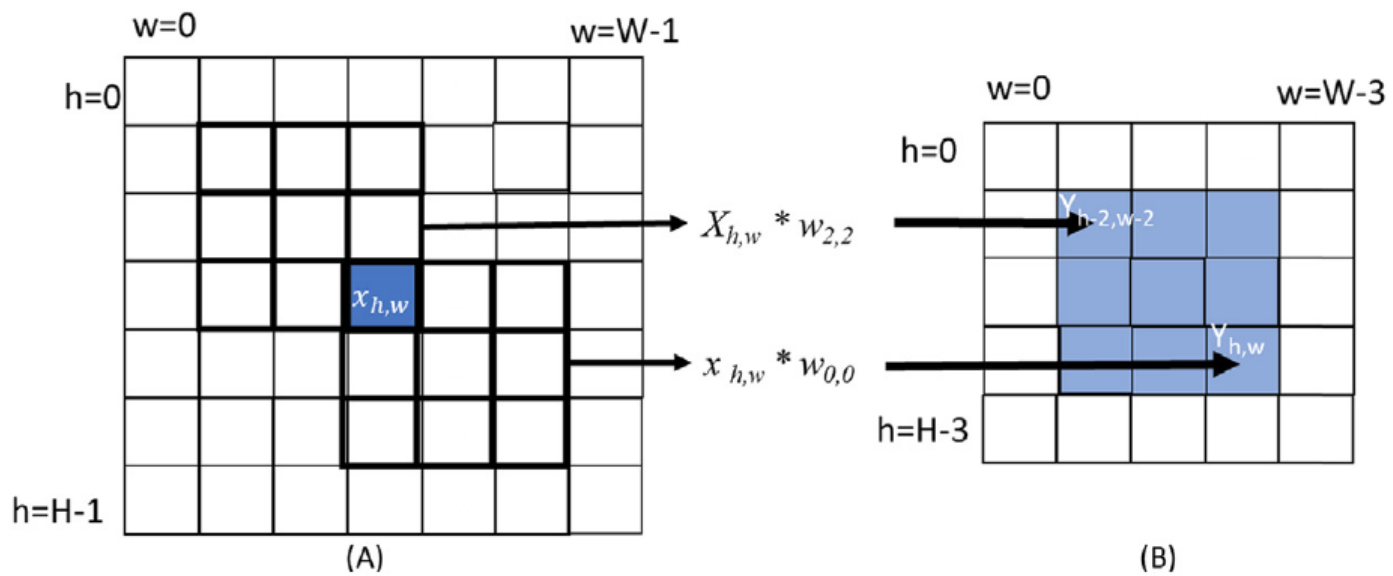


FIGURE 16.9

Convolutional layer. Backpropagation of (A) $\partial E / \partial w$ and (B) $\partial E / \partial x$.

```

01 void convLayer_backward_x_grad(int M, int C, int H_in, int W_in, int K,
                                float* dE_dY, float* W, float* dE_dX) {
02     int H_out = H_in - K + 1;
03     int W_out = W_in - K + 1;
04     for(int c = 0; c < C; c++)
05         for(int h = 0; h < H_in; h++)
06             for(int w = 0; w < W_in; w++)
07                 dE_dX[c, h, w] = 0;

08     for(int m = 0; m < M; m++)
09         for(int h = 0; h < H-1; h++)
10             for(int w = 0; w < W-1; w++)
11                 for(int c = 0; c < C; c++)
12                     for(int p = 0; p < K; p++)
13                         for(int q = 0; q < K; q++)
14                             if(h-p >= 0 && w-p >= 0 && h-p < H_out and w-p < W_OUT)
15                                 dE_dX[c, h, w] += dE_dY[m, h-p, w-p] * W[m, c, k-p, k-q];
16 }

```

FIGURE 16.10

$\frac{\partial E}{\partial x}$ calculation of the backward path of a convolutional layer.

```

01 void convLayer_backward_w_grad(int M, int C, int H, int W, int K, float*
                                dE_dY, float* X, float* dE_dW) {
02     int H_out = H - K + 1;
03     int W_out = W - K + 1;
04     for(int m = 0; m < M; m++)
05         for(int c = 0; c < C; c++)
06             for(int p = 0; p < K; p++)
07                 for(int q = 0; q < K; q++)
08                     dE_dW[m, c, p, q] = 0.;
09
10     for(int m = 0; m < M; m++)
11         for(int h = 0; h < H_out; h++)
12             for(int w = 0; w < W_out; w++)
13                 for(int c = 0; c < C; c++)
14                     for(int p = 0; p < K; p++)
15                         for(int q = 0; q < K; q++)
16                             dE_dW[m, c, p, q] += X[c, h+p, w+q] * dE_dY[m, c, h, w];
17 }

```

FIGURE 16.11

$\frac{\partial E}{\partial w}$ calculation of the backward path of a convolutional layer.

```

01 void convLayer_batched(int N, int M, int C, int H, int W, int K, float* X,
                        float* W, float* Y) {
02     int H_out = H - K + 1;
03     int W_out = W - K + 1;
04     for(int n = 0; n < N; n++)          // for each sample in the mini-batch
05         for(int m = 0; m < M; m++)      // for each output feature map
06             for(int h = 0; h < H_out; h++) // for each output element
07                 for(int w = 0; w < W_out; w++) {
08                     Y[n, m, h, w] = 0;
09                     for (int c = 0; c < C; c++ // sum over all input feature maps
10                         for (int p = 0; p < K; p++)          // KxK filter
11                             for (int q = 0; q < K; q++)
12                                 Y[n,m,h,w] = Y[n,m,h,w] + X[n, c, h+p, w+q]*W[m,c,p,q];
13                 }
14     }

```

FIGURE 16.12

Forward path of a convolutional layer with minibatch training.

```
01  # define TILE_WIDTH 16
02  W_grid = W_out/TILE_WIDTH;  // number of horizontal tiles per output map
03  H_grid = H_out/TILE_WIDTH;  // number of vertical tiles per output map
04  T = H_grid * W_grid;
05  dim3 blockDim(TILE_WIDTH, TILE_WIDTH, 1);
06  dim3 gridDim(M, T, N);
07  ConvLayerForward_Kernel<<< gridDim, blockDim>>>(...);
```

FIGURE 16.13

Host code for launching a convolutional layer kernel.

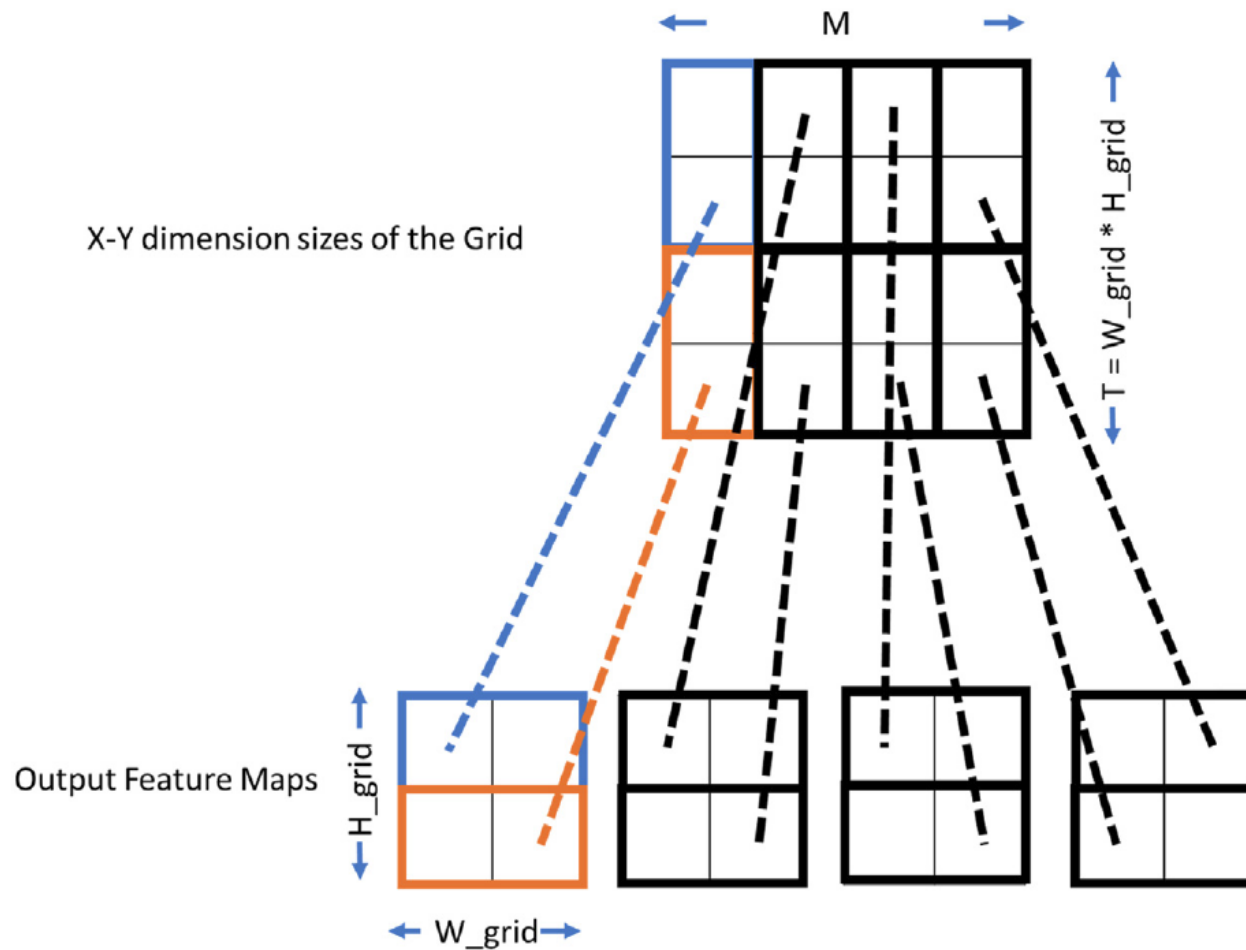


FIGURE 16.14

Mapping output feature map tiles to blocks in the X-Y dimension of the grid.

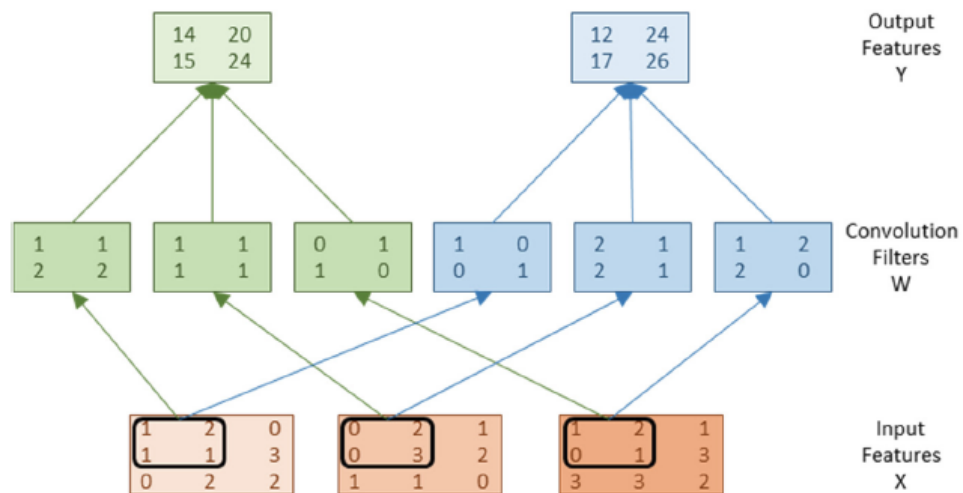
```

01  __global__ void
02  ConvLayerForward_Kernel(int C, int W_grid, int K, float* X, float* W,
                           float* Y) {
03      int m = blockIdx.x;
04      int h = (blockIdx.y / W_grid)*TILE_WIDTH + threadIdx.y;
05      int w = (blockIdx.y % W_grid)*TILE_WIDTH + threadIdx.x;
06      int n = blockIdx.z;
07      float acc = 0.;
08      for (int c = 0; c < C; c++) {          // sum over all input channels
09          for (int p = 0; p < K; p++)        // loop over KxK filter
10              for (int q = 0; q < K; q++)
11                  acc += X[n, c, h + p, w + q] * W[m, c, p, q];
12      }
13      Y[n, m, h, w] = acc;
14  }

```

FIGURE 16.15

Kernel for the forward path of a convolutional layer.



$$\begin{array}{|c|c|c|c|} \hline 1 & 1 & 2 & 2 \\ \hline 1 & 0 & 0 & 1 \\ \hline \end{array}
 \begin{array}{|c|c|c|c|} \hline 1 & 1 & 1 & 1 \\ \hline 2 & 1 & 2 & 1 \\ \hline \end{array}
 \begin{array}{|c|c|c|c|} \hline 0 & 1 & 1 & 0 \\ \hline 1 & 2 & 2 & 0 \\ \hline \end{array}
 *
 \begin{array}{|c|c|c|c|} \hline 1 & 2 & 1 & 1 \\ \hline 2 & 0 & 1 & 3 \\ \hline 1 & 1 & 0 & 2 \\ \hline 1 & 3 & 2 & 2 \\ \hline 0 & 2 & 0 & 3 \\ \hline 2 & 1 & 3 & 2 \\ \hline 0 & 3 & 1 & 1 \\ \hline 3 & 2 & 1 & 0 \\ \hline 1 & 2 & 0 & 1 \\ \hline 2 & 1 & 1 & 3 \\ \hline 0 & 1 & 3 & 3 \\ \hline 1 & 3 & 3 & 2 \\ \hline \end{array}
 =
 \begin{array}{|c|c|c|c|} \hline 14 & 20 & 15 & 24 \\ \hline 12 & 24 & 17 & 26 \\ \hline \end{array}$$

Convolution Filters W Input Features X_unrolled Output Features Y

FIGURE 16.16

Formulation of convolutional layer as GEMM.

```

01 void unroll(int C, int H, int W, int K, float* X, float* X_unroll) {
02     int H_out = H - K + 1;
03     int W_out = W - K + 1;
04     for(int c = 0; c < C; c++) {
05         // Beginning row index of the section for channel C input feature
06         // map in the unrolled matrix
07         w_base = c * (K*K);
08         for(int p = 0; p < K; p++) {
09             for(int q = 0; q < K; q++) {
10                 for(int h = 0; h < H_out; h++) {
11                     int h_unroll = w_base + p*K + q;
12                     for(int w = 0; w < W_out; w++) {
13                         int w_unroll = h * W_out + w;
14                         X_unroll[h_unroll, w_unroll] = X(c, h + p, w + q);
15                     }
16                 }
17             }
18         }
19     }
20 }

```

FIGURE 16.17

A C function that generates the unrolled X matrix. The array accesses are in multidimensional indexing form for clarity and need to be linearized for the code to be compilable.

```

01  __global__ void
02  unroll_Kernel(int C, int H, int W, int K, float* X, float* X_unroll) {
03      int t = blockIdx.x * blockDim.x + threadIdx.x;
04      int H_out = H - K + 1;
05      int W_out = W - K + 1;
06      // Width of the unrolled input feature matrix
07      int W_unroll = H_out * W_out;
08      if (t < C * W_unroll) {
09          // Channel of the input feature map being collected by the thread
10          int c = t / W_unroll;
11          // Column index of the unrolled matrix to write a strip of
12          // input elements into (also, the linearized index of the output
13          // element for which the thread is collecting input elements)
14          int w_unroll = t % W_unroll;
15          // Horizontal and vertical indices of the output element
16          int h_out = w_unroll / W_out;
17          int w_out = w_unroll % W_out;
18          // Starting row index for the unrolled matrix section for channel c
19          int w_base = c * K * K;
20          for(int p = 0; p < K; p++)
21              for(int q = 0; q < K; q++) {
22                  // Row index of the unrolled matrix for the thread to write
23                  // the input element into for the current iteration
24                  int h_unroll = w_base + p*K + q;
25                  X_unroll[h_unroll, w_unroll] = X[c, h_out + p, w_out + q];
26              }
27      }
28  }

```

FIGURE 16.18

A CUDA kernel implementation for unrolling input feature maps. The array accesses are in multidimensional indexing form for clarity and need to be linearized for the code to be compilable.

Table 16.1 Convolution parameters for CUDNN. Note that the CUDNN naming convention is slightly different from what we used in previous sections.

Parameter	Meaning
N	Number of images in minibatch
C	Number of input feature maps
H	Height of input image
W	Width of input image
K	Number of output feature maps
R	Height of filter
S	Width of filter
u	Vertical stride
v	Horizontal stride
pad_h	Height of zero padding
pad_w	Width of zero padding