

Mathematics of AI

Compressive sensing, convolutions, other random interesting stuff

Goodfellow et al, *Deep Learning*

Random papers, textbooks, and online resources

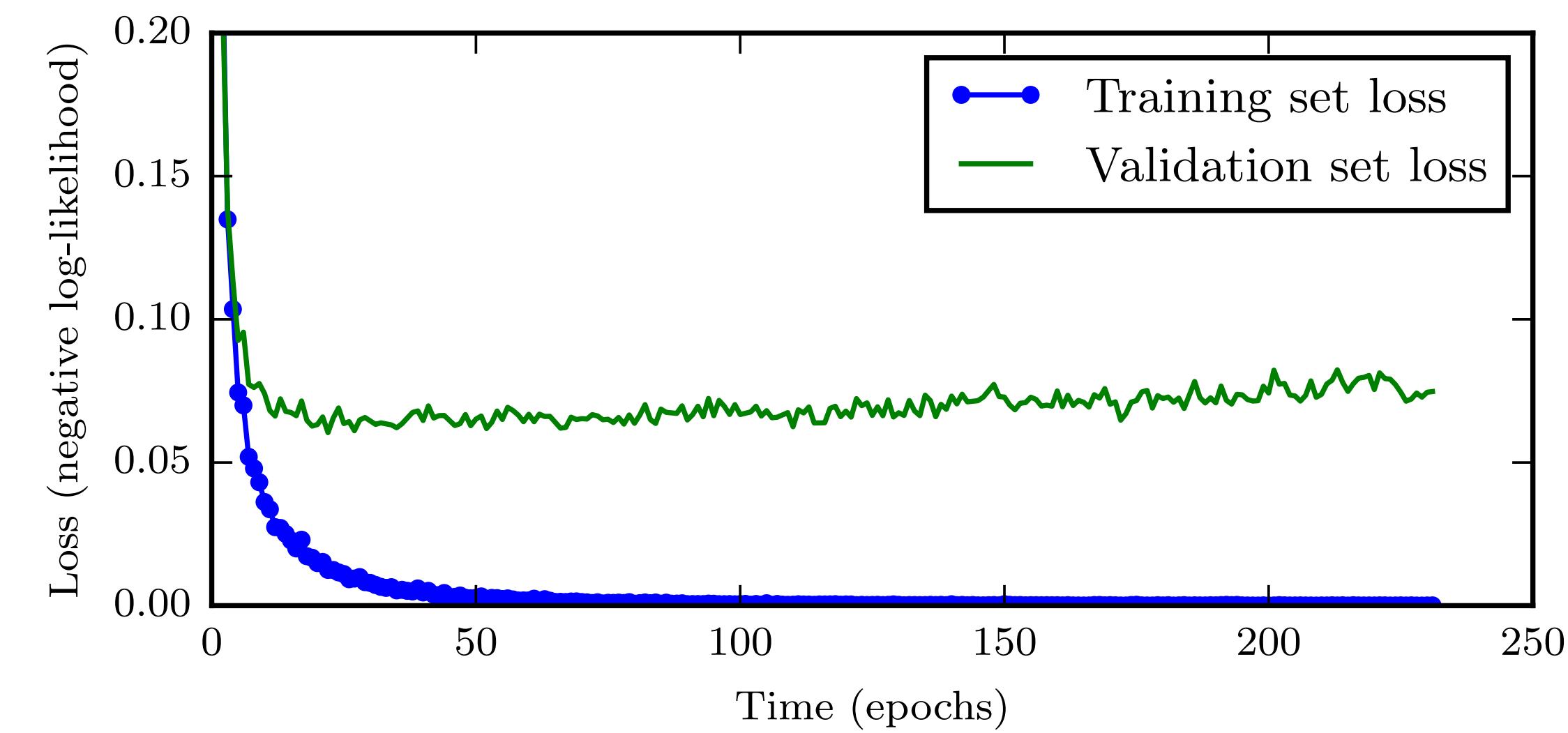
Definition

- “Regularization is any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.”

Regularisation

In deep learning

- “Just get more data” – Guy from Google
- Early stopping, i.e., “cheating”
- ... (bagging, multi-task learning, ...)
- Sparse representations! (cf. compressive sensing...)
- Dropout



Sparse Representations

$$\begin{bmatrix} 18 \\ 5 \\ 15 \\ -9 \\ -3 \end{bmatrix} = \begin{bmatrix} 4 & 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & -1 & 0 & 3 & 0 \\ 0 & 5 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & -4 \\ 1 & 0 & 0 & 0 & -5 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ -2 \\ -5 \\ 1 \\ 4 \end{bmatrix} \quad (7.46)$$

$\mathbf{y} \in \mathbb{R}^m \qquad \qquad \mathbf{A} \in \mathbb{R}^{m \times n} \qquad \qquad \mathbf{x} \in \mathbb{R}^n$

Sparse Representations

$$\begin{bmatrix} -14 \\ 1 \\ 19 \\ 2 \\ 23 \end{bmatrix} = \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix} \quad (7.47)$$

$\boldsymbol{y} \in \mathbb{R}^m \qquad \boldsymbol{B} \in \mathbb{R}^{m \times n} \qquad \boldsymbol{h} \in \mathbb{R}^n$

Sparse Representations

$$\begin{bmatrix} -14 \\ 1 \\ 19 \\ 2 \\ 23 \end{bmatrix} = \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix} \quad (7.47)$$

$\mathbf{y} \in \mathbb{R}^m \qquad \qquad \mathbf{B} \in \mathbb{R}^{m \times n} \qquad \qquad \mathbf{h} \in \mathbb{R}^n$

Sparser representation of
data ==> greater
generalisability

$$\arg \min_{\mathbf{h}, \|\mathbf{h}\|_0 < k} \|\mathbf{x} - \mathbf{W}\mathbf{h}\|^2,$$

Sparse Representations

$$\begin{bmatrix} -14 \\ 1 \\ 19 \\ 2 \\ 23 \end{bmatrix} = \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix} \quad (7.47)$$

$\mathbf{y} \in \mathbb{R}^m \qquad \mathbf{B} \in \mathbb{R}^{m \times n} \qquad \mathbf{h} \in \mathbb{R}^n$

Learn:

- Dictionary \mathbf{W} (or \mathbf{D})
- Representation \mathbf{h} (or \mathbf{f})

$$\arg \min_{\mathbf{h}, \|\mathbf{h}\|_0 < k} \|\mathbf{x} - \mathbf{W}\mathbf{h}\|^2,$$

Sparse Representations

The Importance of Encoding Versus Training with Sparse Coding and Vector Quantization

Adam Coates

Andrew Y. Ng

Stanford University, 353 Serra Mall, Stanford, CA 94305

ACOATES@CS.STANFORD.EDU

ANG@CS.STANFORD.EDU

| TRAIN / ENCODER | TEST ACC. |
|-------------------------------------|--------------|
| RP / T | 79.1% |
| SC / SC | 78.8% |
| SC / T | 78.9% |
| OMP-1 / SC | 78.8% |
| OMP-1 / T | 79.4% |
| OMP-10 / T | 80.1% |
| OMP-1 / T ($d = 6000$) | 81.5% |
| (COATES ET AL., 2011) 1600 FEATURES | 77.9% |
| (COATES ET AL., 2011) 4000 FEATURES | 79.6% |
| IMPROVED LCC (YU & ZHANG, 2010) | 74.5% |
| CONV. DBN (Krizhevsky, 2010) | 78.9% |
| DEEP NN (Ciresan et al., 2011) | 80.49% |

Lots of ways to learn
dictionaries and
representations!

<https://openreview.net/forum?id=By4af0WO-B>

(Goodfellow 2016)

results using dictionaries created from random noise, randomly sampled exemplars, and vector quantization and show that the last two yield perfectly usable dictionaries in every case.

3. Learning framework

Given an unsupervised learning algorithm, we learn a new feature representation from unlabeled data by employing a common framework. Our feature learning framework is like the patch-based system presented in (Coates et al., 2011) with a few modifications. It shares most of its key components with prior work in visual word models (Csurka et al., 2004; Lazebnik et al., 2006; Agarwal & Triggs, 2006).

In order to generate a set of features, our system first accumulates a batch of small image patches or image descriptors harvested from unlabeled data. When learning from raw pixels, we extract 6 pixel square patches, yielding a bank of vectors that are then normalized¹ and ZCA whitened (Hyvarinen & Oja, 2000) (retaining full variance). If we are learning from SIFT descriptors, we simply take single descriptors to form a bank of 128-dimensional vectors.

Given the batch of input vectors, $x^{(i)} \in \mathbb{R}^n$, an unsupervised learning algorithm is then applied to learn a dictionary of d elements, $D \in \mathbb{R}^{n \times d}$, where each column $D^{(j)}$ is one element. In order to make all of our algorithms consistent (so that we may freely change the choice of encoder), we will make certain that each of the algorithms we use produces normalized dictionary elements: $\|D^{(j)}\|_2^2 = 1$.

In this work, we will use the following unsupervised learning algorithms for training the dictionary D :

1. **Sparse coding (SC):** We train the dictionary using the L1-penalized sparse coding formulation. That is, we optimize

$$\min_{D,s(i)} \sum_i \|Ds^{(i)} - x^{(i)}\|_2^2 + \lambda \|s^{(i)}\|_1 \quad (1)$$

subject to $\|D^{(j)}\|_2^2 = 1, \forall j$

using alternating minimization over the sparse codes, $s^{(i)}$, and the dictionary, D . We use the coordinate descent algorithm to solve for the sparse codes (Wu & Lange, 2008).

2. **Orthogonal matching pursuit (OMP-k):** Similar to sparse coding, the dictionary is trained

using an alternating minimization of

$$\min_{D,s(i)} \sum_i \|Ds^{(i)} - x^{(i)}\|_2^2 \quad (2)$$

subject to $\|D^{(j)}\|_2^2 = 1, \forall j$

and $\|s^{(i)}\|_0 \leq k, \forall i$

where $\|s^{(i)}\|_0$ is the number of non-zero elements in $s^{(i)}$. In this case, the codes $s^{(i)}$ are computed (approximately) using Orthogonal Matching Pursuit (Pati et al., 1993; Blumensath & Davies, 2007) to compute codes with at most k non-zeros (which we refer to as “OMP-k”). For a single input $x^{(i)}$, OMP-k begins with $s^{(i)} = 0$ and at each iteration greedily selects an element of $s^{(i)}$ to be made non-zero to minimize the residual reconstruction error. After each selection, $s^{(i)}$ is updated to minimize $\|Ds^{(i)} - x^{(i)}\|_2^2$ over $s^{(i)}$ allowing only the selected elements to be non-zero.

Importantly, OMP-1 is a form of “gain-shape vector quantization” (and is similar to K-means when the data $x^{(i)}$ and the dictionary elements $D^{(j)}$ all have unit length). Specifically, it chooses $k = \arg \max_j |D^{(j)\top} x^{(i)}|$, then sets $s_k^{(i)} = D^{(k)\top} x^{(i)}$ and all other elements of $s^{(i)}$ to 0. Holding these “one hot” codes fixed, it is then easy to solve for the optimal D in (2).

3. Sparse RBMs and sparse auto-encoders (RBM, SAE):

In some of our experiments, we train sparse RBMs (Hinton et al., 2006) and sparse auto-encoders (Ranzato et al., 2007; Bengio et al., 2006), both using a logistic sigmoid non-linearity $g(Wx + b)$. These algorithms yield a set of weights W and biases b . To obtain the dictionary, D , we simply discard the biases and take $D = W^\top$, then normalize the columns of D .

4. Randomly sampled patches (RP):

We also use a heuristic method for populating the dictionary: we fill the columns of D with normalized vectors sampled randomly from amongst the $x^{(i)}$.

5. Random weights (R):

It has also been shown that completely random weights can perform surprisingly well in some tasks (Jarrett et al., 2009; Saxe et al., 2010). Thus, we have also tried filling the columns of D with vectors sampled from a unit normal distribution (subsequently normalized to unit length).

After running any of the above training procedures, we have a dictionary D . We must then define an “encoder” that, given D , maps a new input vector x to

¹We subtract the mean and divide by the standard deviation of the pixel values.

algorithms consistent (so that we may freely change the choice of encoder), we will make certain that each of the algorithms we use produces normalized dictionary elements: $\|D^{(j)}\|_2^2 = 1$.

In this work, we will use the following unsupervised learning algorithms for training the dictionary D :

1. **Sparse coding (SC):** We train the dictionary using the L1-penalized sparse coding formulation.

That is, we optimize

$$\min_{D, s^{(i)}} \sum_i \|Ds^{(i)} - x^{(i)}\|_2^2 + \lambda \|s^{(i)}\|_1 \quad (1)$$

subject to $\|D^{(j)}\|_2^2 = 1, \forall j$

constrained
optimisation

using alternating minimization over the sparse codes, $s^{(i)}$, and the dictionary, D . We use the coordinate descent algorithm to solve for the sparse codes (Wu & Lange, 2008).

2. **Orthogonal matching pursuit (OMP-k):** Similar to sparse coding, the dictionary is trained

¹We subtract the mean and divide by the standard deviation of the pixel values.

3. **Sparse RBMs and sparse auto-encoders (RBM, SAE):** In some of our experiments, we train sparse RBMs (Hinton et al., 2006) and sparse auto-encoders (Ranzato et al., 2007; Bengio et al., 2006), both using a logistic sigmoid non-linearity $g(Wx + b)$. These algorithms yield a set of weights W and biases b . To obtain the dictionary, D , we simply discard the biases and take $D = W^\top$, then normalize the columns of D .

4. **Randomly sampled patches (RP):** We also use a heuristic method for populating the dictionary: we fill the columns of D with normalized vectors sampled randomly from amongst the $x^{(i)}$.

5. **Random weights (R):** It has also been shown that completely random weights can perform surprisingly well in some tasks (Jarrett et al., 2009; Saxe et al., 2010). Thus, we have also tried filling the columns of D with vectors sampled from a unit normal distribution (subsequently normalized to unit length).

After running any of the above training procedures, we have a dictionary D . We must then define an “encoder” that, given D , maps a new input vector x to

The Importance of Encoding Versus Training with Sparse Coding and Vector Quantization

results using dictionaries created from random noise, randomly sampled exemplars, and vector quantization and show that the last two yield perfectly usable dictionaries in every case.

constrained
optimisation

3. Learning framework

Given an unsupervised learning algorithm, we learn a new feature representation from unlabeled data by employing a common framework. Our feature learning framework is like the patch-based system presented in (Coates et al., 2011) with a few modifications. It shares most of its key components with prior work in visual word models (Csurka et al., 2004; Lazebnik et al., 2006; Agarwal & Triggs, 2006).

In order to generate a set of features, our system first accumulates a batch of small image patches or image descriptors harvested from unlabeled data. When learning from raw pixels, we extract 6 pixel square patches, yielding a bank of vectors that are then normalized¹ and ZCA whitened (Hyvarinen & Oja, 2000)

using an alternating minimization of

$$\min_{D, s^{(i)}} \sum_i \|Ds^{(i)} - x^{(i)}\|_2^2 \quad (2)$$

subject to $\|D^{(j)}\|_2^2 = 1, \forall j$
and $\|s^{(i)}\|_0 \leq k, \forall i$

where $\|s^{(i)}\|_0$ is the number of non-zero elements in $s^{(i)}$. In this case, the codes $s^{(i)}$ are computed (approximately) using Orthogonal Matching Pursuit (Pati et al., 1993; Blumensath & Davies, 2007) to compute codes with at most k non-zeros (which we refer to as “OMP-k”). For a single input $x^{(i)}$, OMP-k begins with $s^{(i)} = 0$ and at each iteration greedily selects an element of $s^{(i)}$ to be made non-zero to minimize the residual reconstruction error. After each selection, $s^{(i)}$ is updated to minimize $\|Ds^{(i)} - x^{(i)}\|_2^2$ over $s^{(i)}$ allowing only the selected elements to be non-zero.

Importantly, OMP-1 is a form of “gain-shape vector quantization” (similar to K-

Or a regularisation problem

Matt's notes

Reconstruction, inverse problems, and regularisation

Transformation is just

$$\mathbf{z} = W\mathbf{x}, \quad (6)$$

but if this is an invertible transform then W is a square, non-singular matrix. But if we are working, for instance, with an overcomplete dictionary, the transform matrix W might be non-square. Also in many problems we are forced into measuring (by cost or other practicalities) such that we measure through some transformation.

Hence, a very common problem is reconstruction of a signal \mathbf{x} from a transformed measurement \mathbf{z} . Classic problems in this form are sometimes called *tomography* problems by analogy to CAT (Computer Axial Tomography) scans. These reconstruction problems are inverse problems (I shall stick to linear problems for now). They have a set of general characteristics:

- The equation above is more realistically $\mathbf{z} = W\mathbf{x} + \epsilon$, where ϵ is noise. Hence, there may not actually be an exact solution to the set of equations.
- The matrix W is usually not square and in many cases the length of \mathbf{z} is much smaller than \mathbf{x} so the equations are (massively) underconstrained.

So we can't solve this using the usual linear algebra you have been taught.

The easiest way to create a framework for all of these problems is to consider them as optimisation problems, i.e., we write them as

$$\operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) \text{ such that } \mathbf{z} = W\mathbf{x}. \quad (7)$$

but for our purposes we solve a new optimisation problem

$$\operatorname{argmin}_{\mathbf{x}} f(\mathbf{x}) + \lambda \|\mathbf{z} - W\mathbf{x}\|. \quad (8)$$

That is we solve a problem that tries to find a good solution that closely satisfies the measurement constraints (the extra term in the objective is a type of Lagrange multiplier).

Note the extra parameter λ . This lets you perform a tradeoff between your "model" implied by $f(\cdot)$ and your measurements. So

- If you think the measurements are very noisy, you might choose a smaller λ .
- If you think the measurements are very accurate, you might choose a larger λ .

The secret of all of these, however, lies in the details of $f(\cdot)$ and the norm we use $\|\cdot\|$. There are very many alternatives.

cf. Compressive sensing

Matt's notes

Sensing and compression as one action = compressive sensings

Now we get to put all of these ideas together.

1. Remember that transformation and sensing can happen as one step in hardware (e.g., JPEG) so
 - I will assume that my signal \mathbf{x} is compressible (after some transformation of the raw signal), and
 - I will also use a random projection matrix A to perform a dimension reduction on the signal to get

$$\mathbf{y} = A\mathbf{x} \tag{9}$$

where $\mathbf{x} \in R^n$ and $\mathbf{y} \in R^k$ where $k \ll n$ (i.e., A is a $k \times n$ matrix).

2. Now I want to reconstruct an approximation $\hat{\mathbf{x}}$ from \mathbf{y} and I shall suggest the following:

$$\operatorname{argmin}_{\mathbf{x}} \lambda \|\mathbf{x}\|_0 + \|\mathbf{z} - W\mathbf{x}\|_2^2. \tag{10}$$

(actually use
an L1 norm)

Compressive sensing

Linkages throughout the course

- CS is a constrained optimisation problem
- ... which itself can be cast as unconstrained optimisation with a regularisation (penalty term)
- CS (with the L1 penalty term) is essentially LASSO (L2 penalty is ridge regression)
- Also essentially the same thing as “sparse representations” —> regularisation for deep learning (and works well)
- “Deep” methods essentially do the same thing as CS (e.g., <https://arxiv.org/abs/1905.06723>)

Convolutions and neural networks

Key concepts

- For images (and audio, and ...), wouldn't DCT etc make sense for layers?
- Stride, padding, and dilation
- Pooling

Convolutions and neural networks

Why

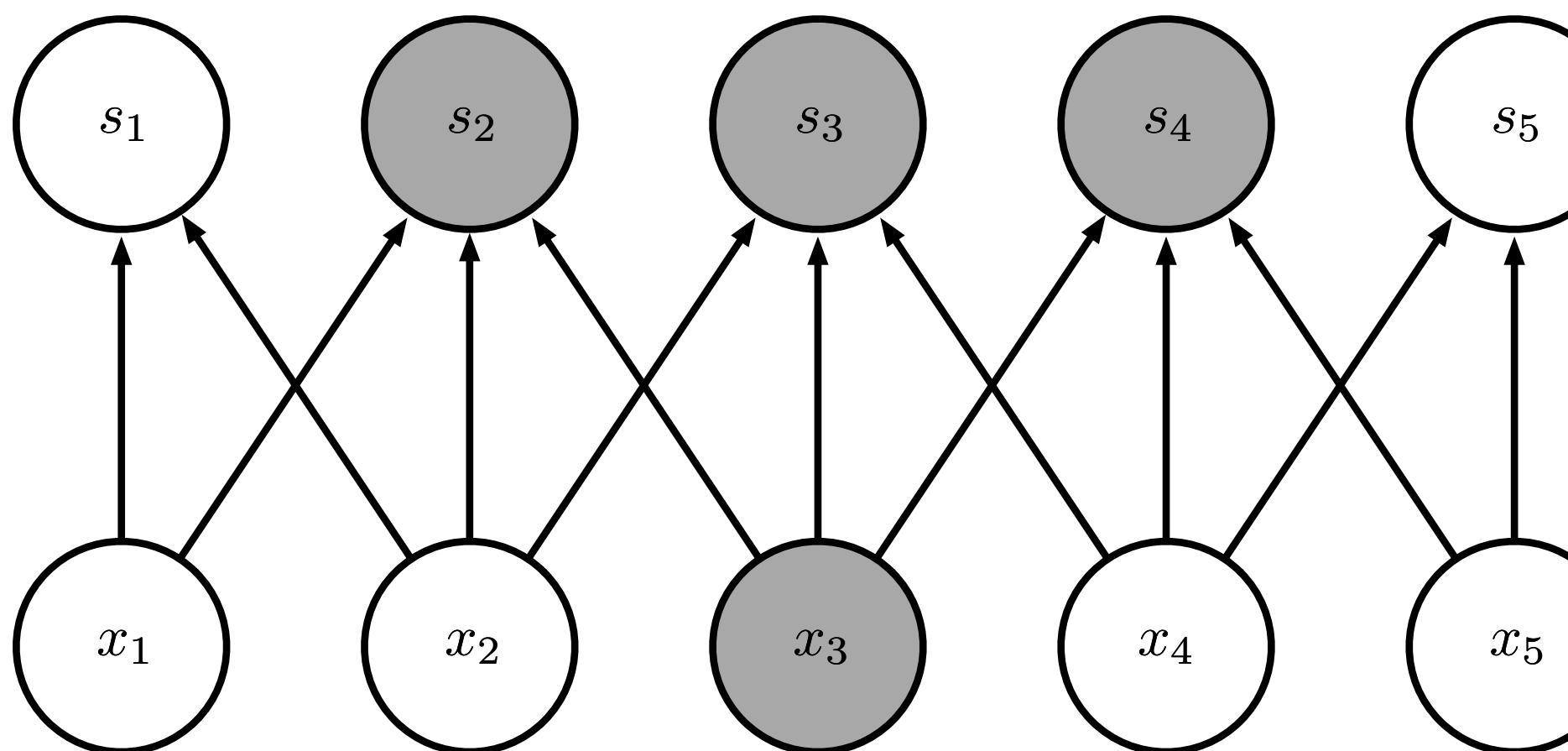
- Consider fully-connected network for speech classification: 2min of 8kHz audio, hidden layers with 10k, 1k neurons

$$\hat{f}(x) = \text{sign}(2\sigma_{\text{sigmoid}}(W_3\sigma_{\text{relu}}(W_2\sigma_{\text{relu}}(W_1x + b_1) + b_2) + b_3) - 1).$$

- Only need to use local data
- (Even worse for images!)

Sparse Connectivity

Sparse
connections
due to small
convolution
kernel



Dense
connections

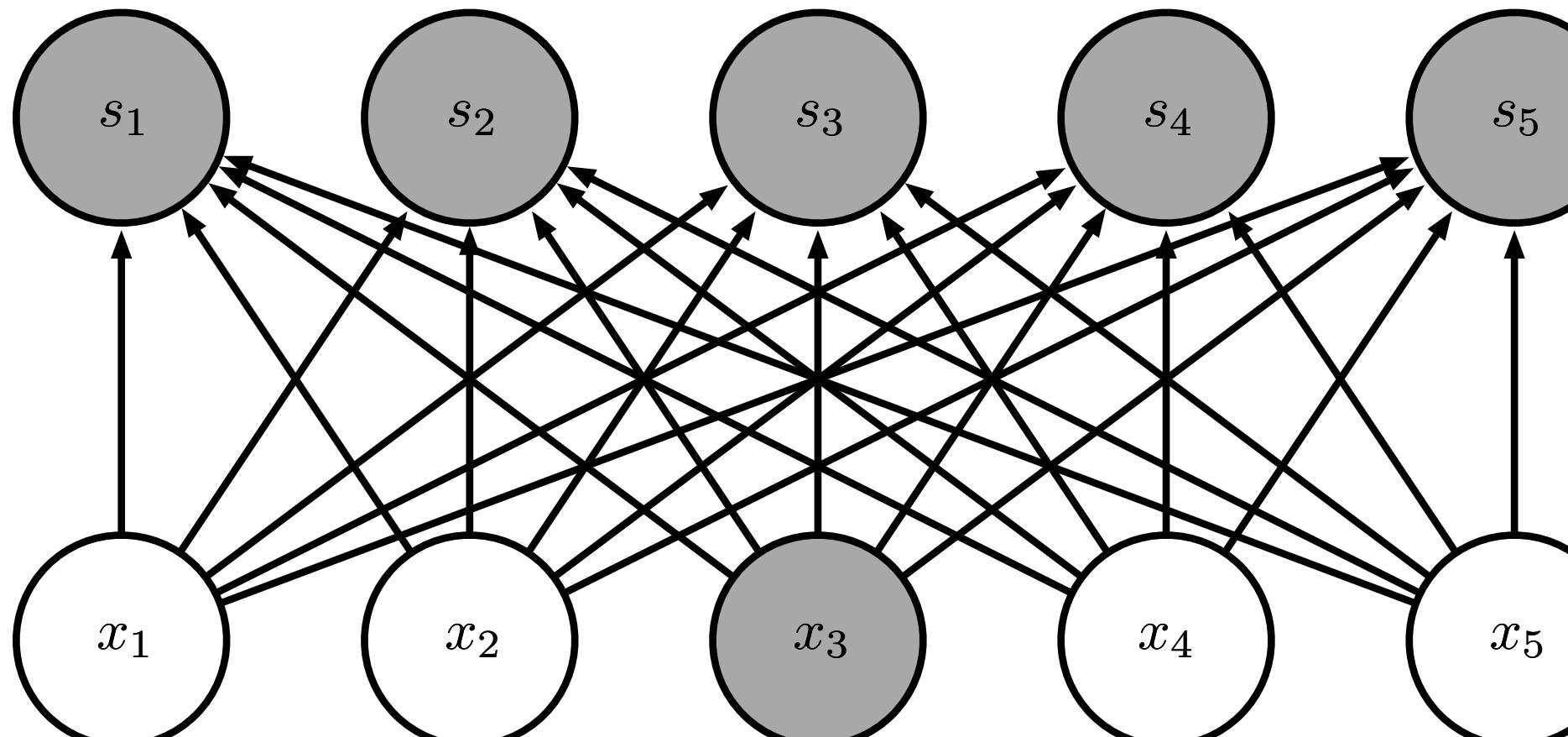
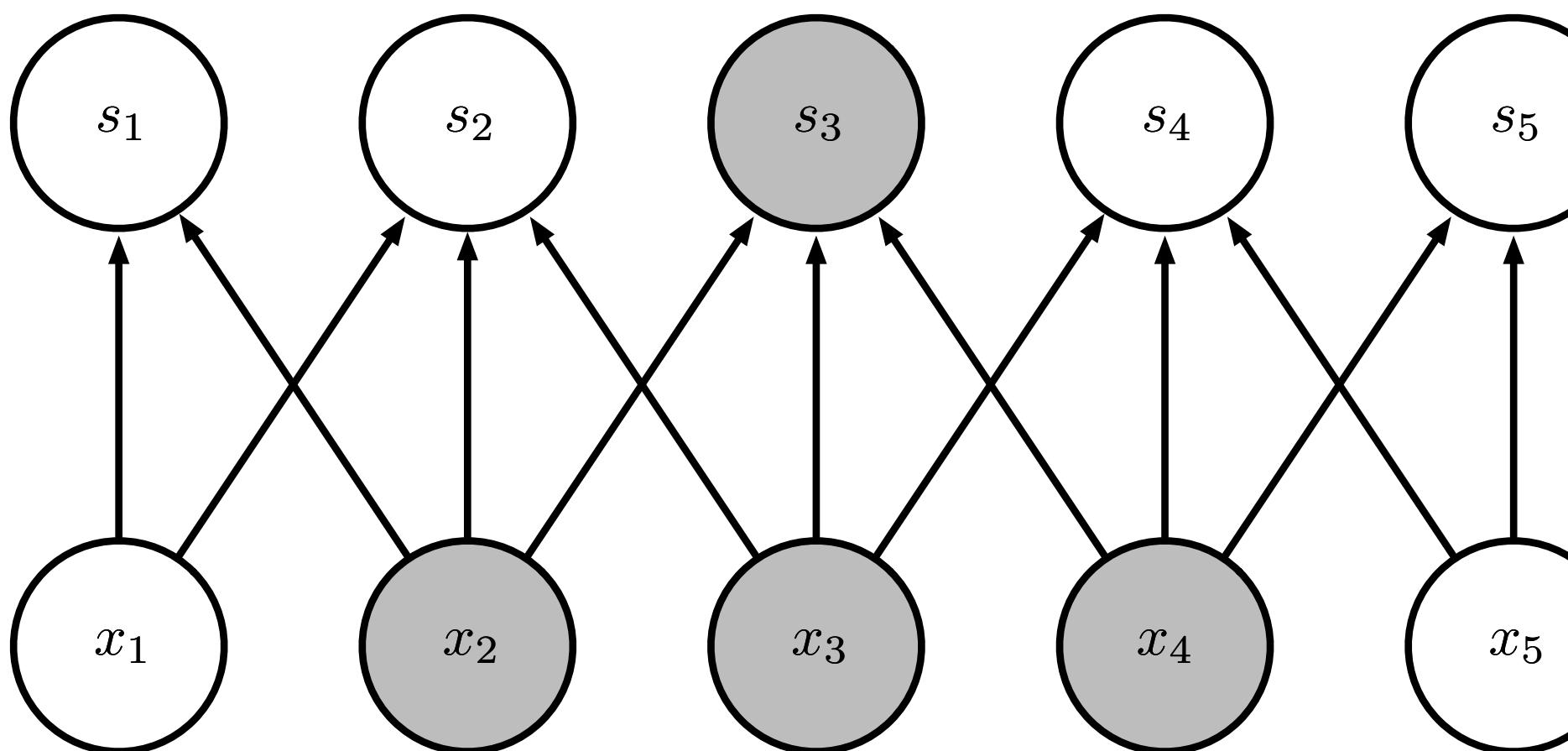


Figure 9.2

Sparse Connectivity

Sparse
connections
due to small
convolution
kernel



Dense
connections

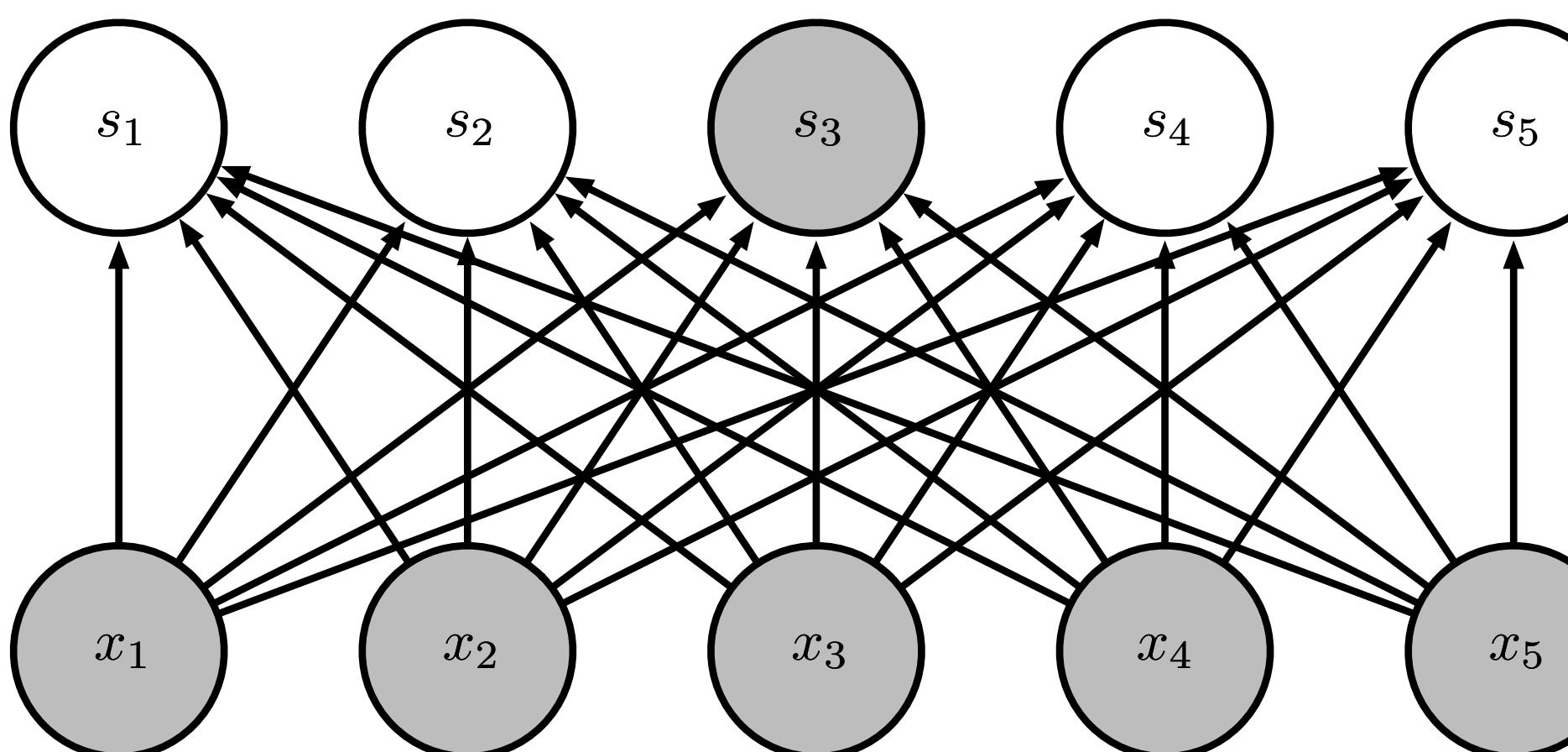


Figure 9.3

2D Convolution

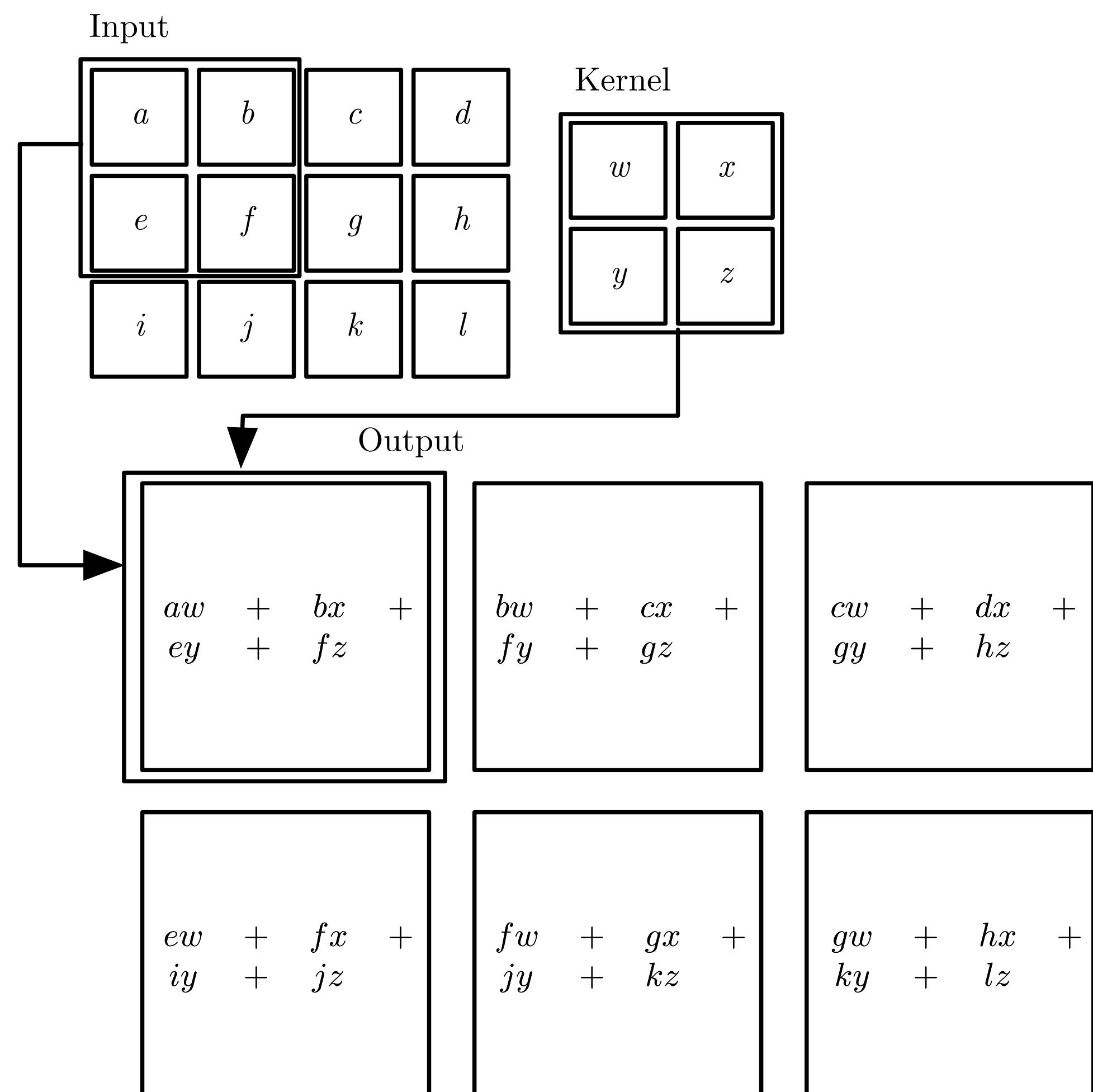


Figure 9.1

(Goodfellow 2016)

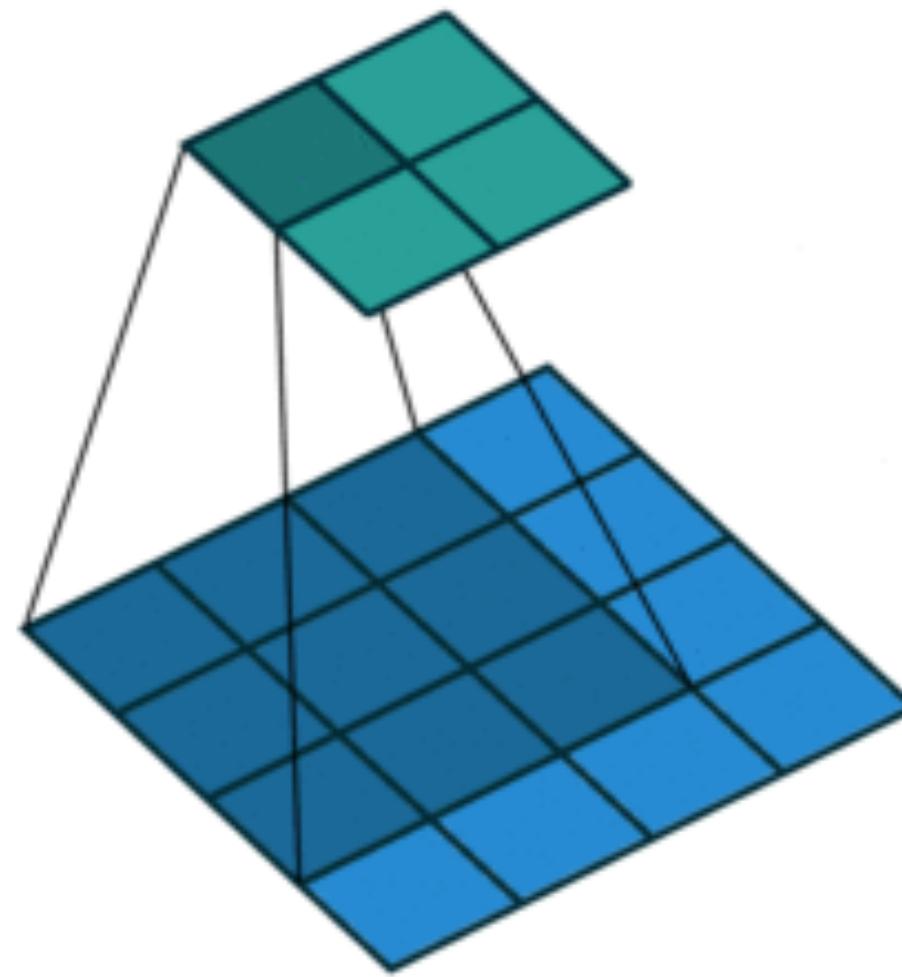
Convolutions in neural networks

Key concepts

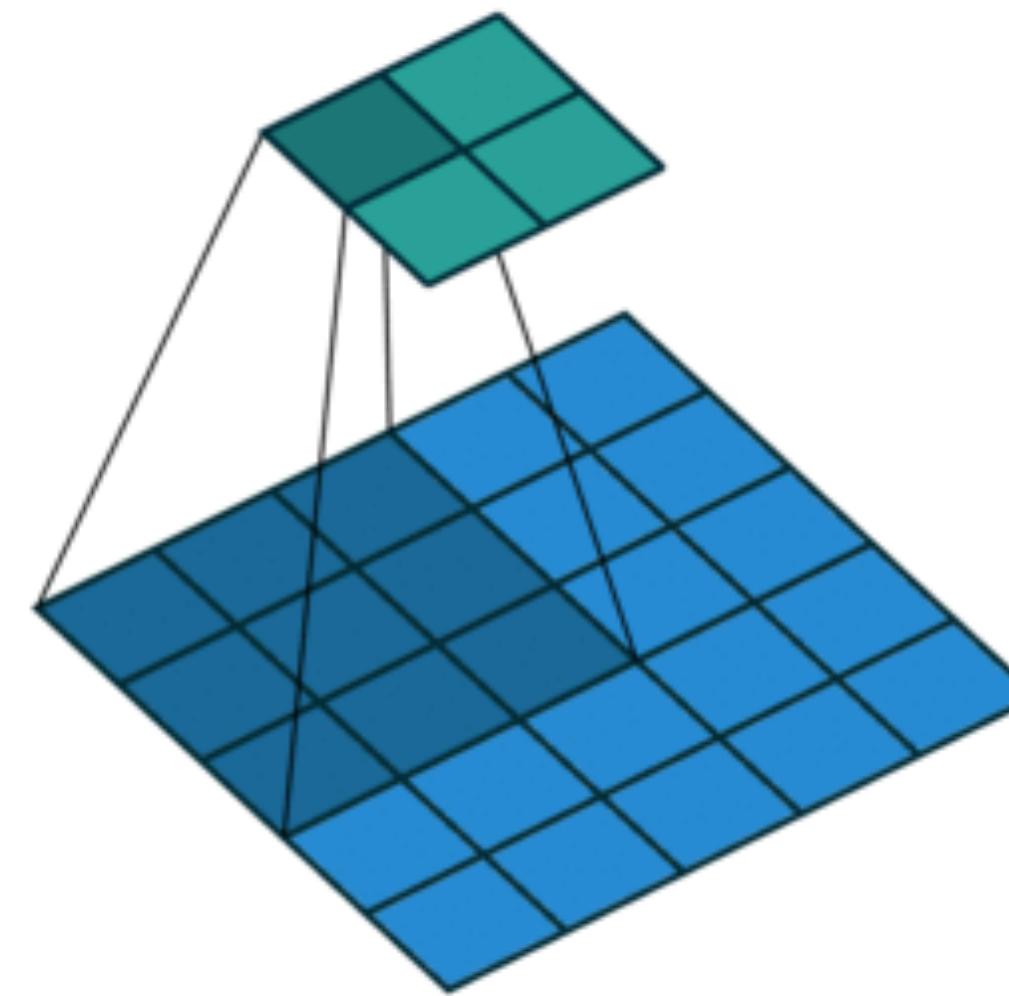
- Strides
- Padding
- Dilation
- Pooling

Convolutions in neural networks

Strides



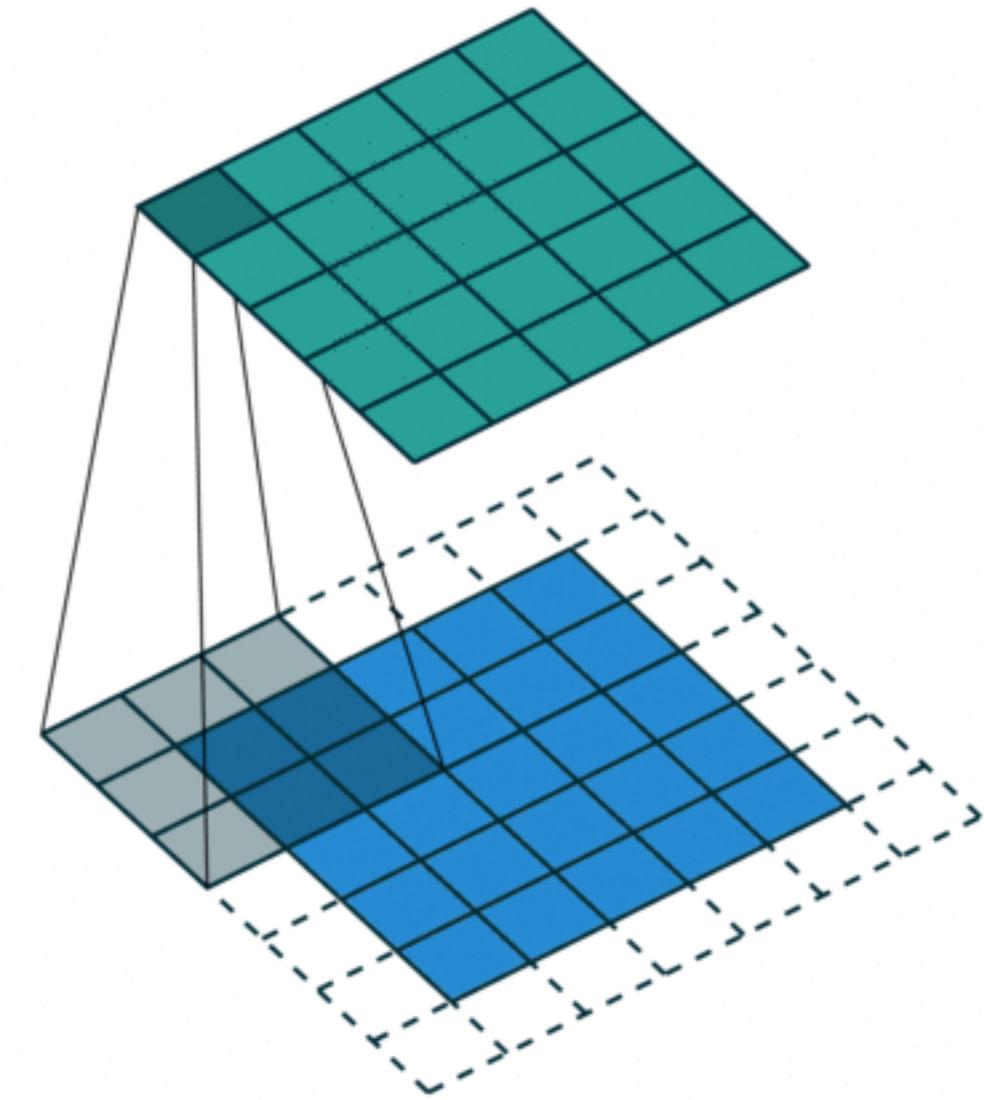
Stride length = 1



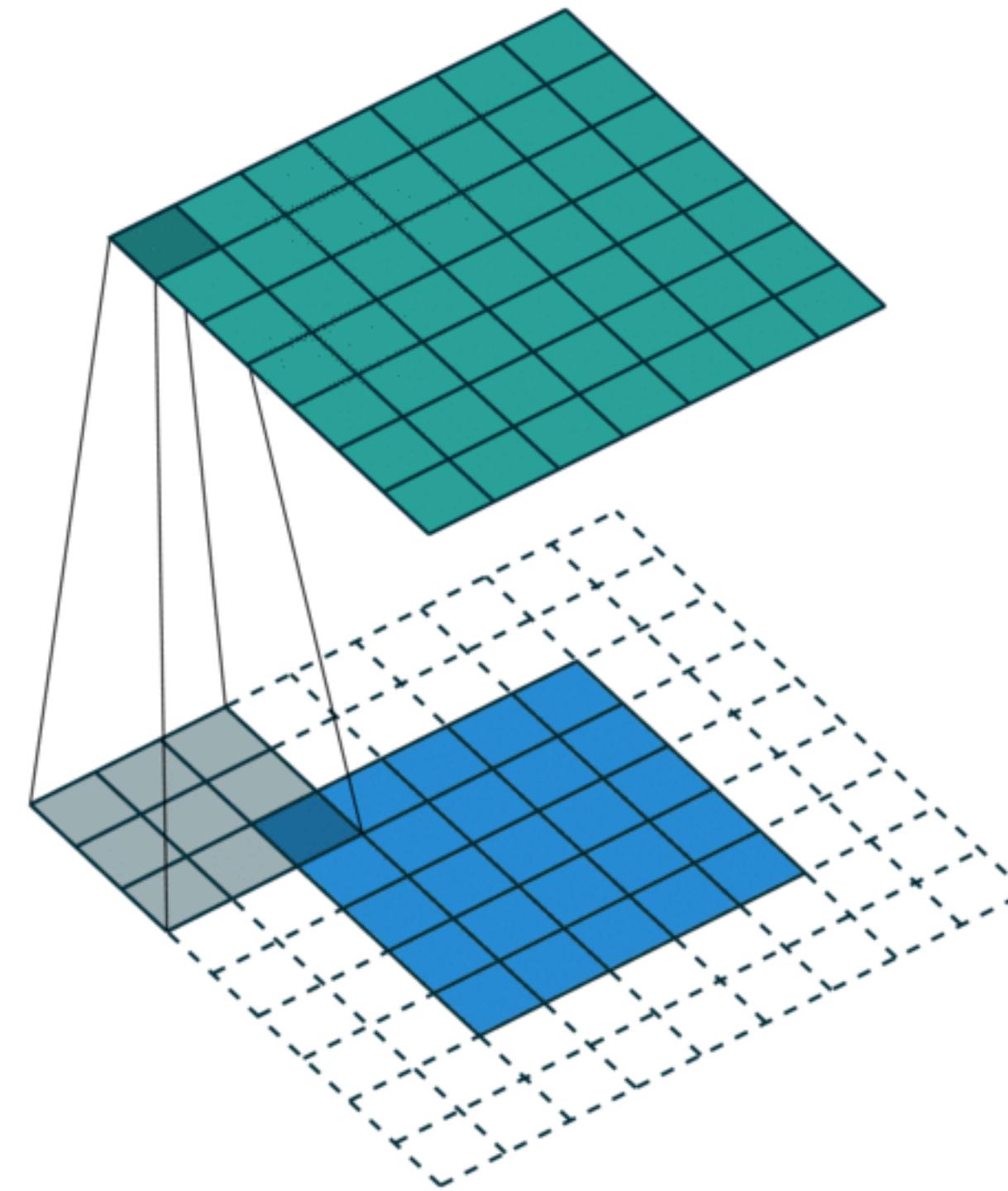
Stride length = 2

Convolutions in neural networks

Padding



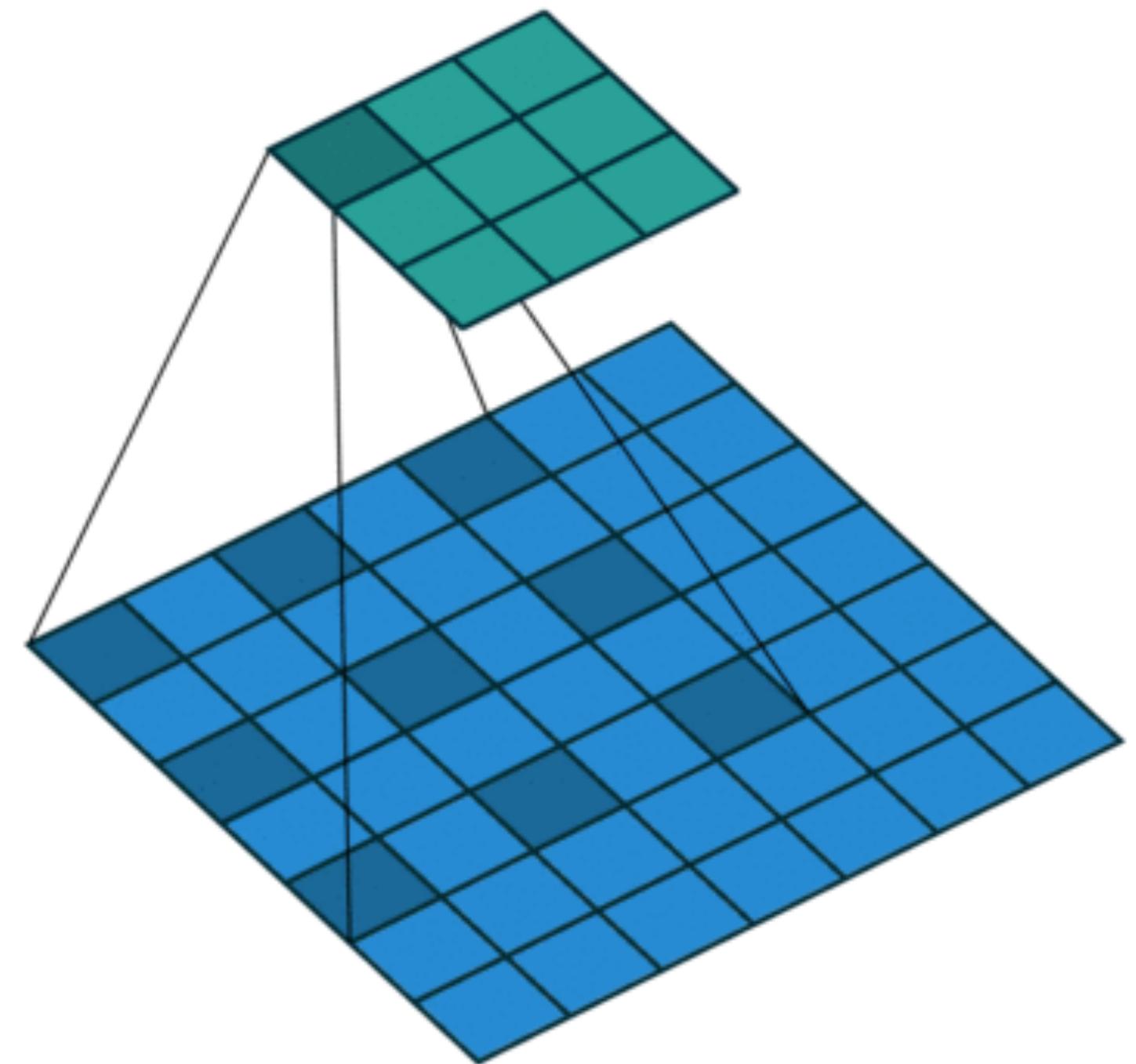
Padding length = 1



Padding length = 2

Convolutions in neural networks

Dilation



Dilation factor = 2

Convolutions in neural networks

Pooling

What do we do to get each pixel in the output?

- Max pooling
- Average pooling
- L2 norm (or your favourite)
- Weighted (e.g., Gaussian)

Convolutions in neural networks

Pooling

| | | |
|---|---|---|
| 0 | 1 | 2 |
| 2 | 2 | 0 |
| 0 | 1 | 2 |

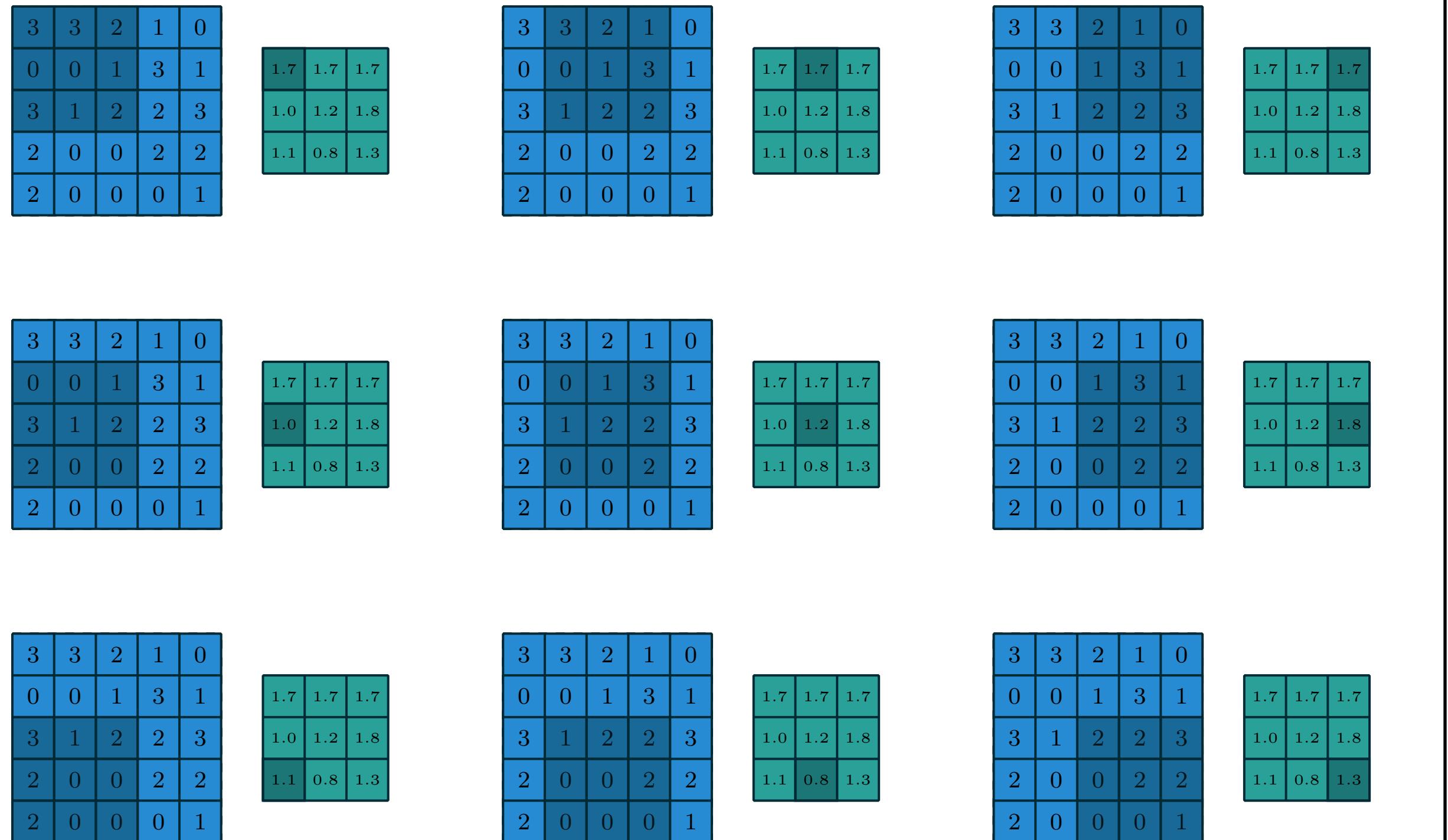


Figure 1.5: Computing the output values of a 3×3 average pooling operation on a 5×5 input using 1×1 strides.

Convolutions in neural networks

Pooling

| | | |
|---|---|---|
| 0 | 1 | 2 |
| 2 | 2 | 0 |
| 0 | 1 | 2 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 1.7 | 1.7 | 1.7 |
| 1.0 | 1.2 | 1.8 |
| 1.1 | 0.8 | 1.3 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 1.7 | 1.7 | 1.7 |
| 1.0 | 1.2 | 1.8 |
| 1.1 | 0.8 | 1.3 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 1.7 | 1.7 | 1.7 |
| 1.0 | 1.2 | 1.8 |
| 1.1 | 0.8 | 1.3 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 1.7 | 1.7 | 1.7 |
| 1.0 | 1.2 | 1.8 |
| 1.1 | 0.8 | 1.3 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 1.7 | 1.7 | 1.7 |
| 1.0 | 1.2 | 1.8 |
| 1.1 | 0.8 | 1.3 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 1.7 | 1.7 | 1.7 |
| 1.0 | 1.2 | 1.8 |
| 1.1 | 0.8 | 1.3 |

Convolutions in neural networks

Pooling

| | | |
|---|---|---|
| 0 | 1 | 2 |
| 2 | 2 | 0 |
| 0 | 1 | 2 |

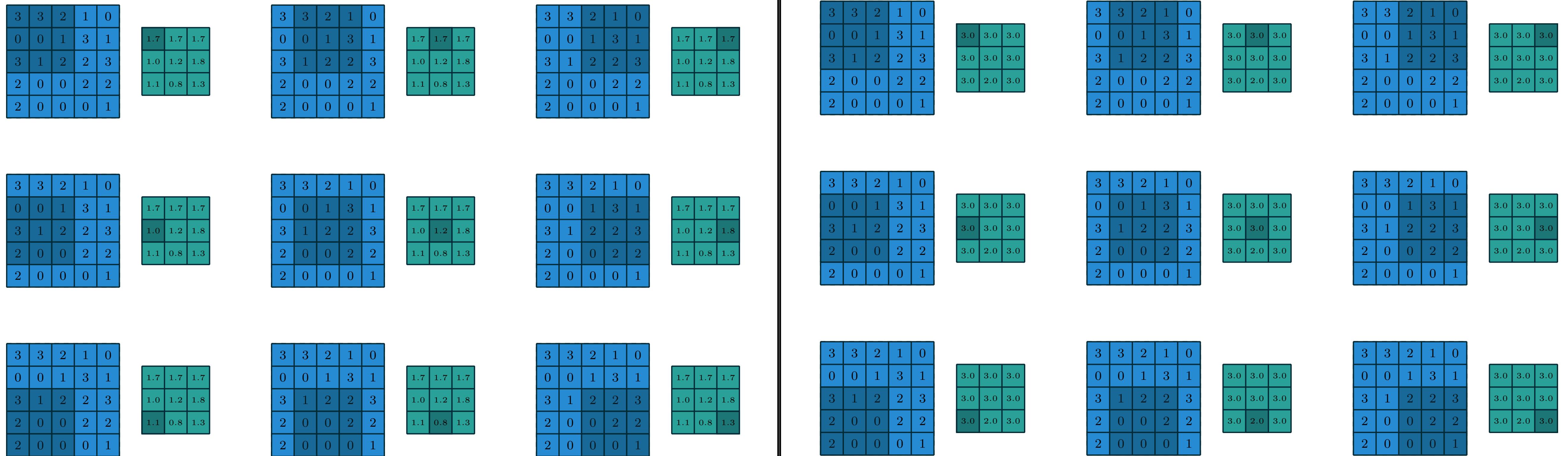


Figure 1.5: Computing the output values of a 3×3 average pooling operation on a 5×5 input using 1×1 strides.

Figure 1.6: Computing the output values of a 3×3 max pooling operation on a 5×5 input using 1×1 strides.

Convolutions in neural networks

Pooling

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 3.0 | 3.0 | 3.0 |
| 3.0 | 3.0 | 3.0 |
| 3.0 | 2.0 | 3.0 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 3.0 | 3.0 | 3.0 |
| 3.0 | 3.0 | 3.0 |
| 3.0 | 2.0 | 3.0 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 3.0 | 3.0 | 3.0 |
| 3.0 | 3.0 | 3.0 |
| 3.0 | 2.0 | 3.0 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 3.0 | 3.0 | 3.0 |
| 3.0 | 3.0 | 3.0 |
| 3.0 | 2.0 | 3.0 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 3.0 | 3.0 | 3.0 |
| 3.0 | 3.0 | 3.0 |
| 3.0 | 2.0 | 3.0 |

| | | | | |
|---|---|---|---|---|
| 3 | 3 | 2 | 1 | 0 |
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

| | | |
|-----|-----|-----|
| 3.0 | 3.0 | 3.0 |
| 3.0 | 3.0 | 3.0 |
| 3.0 | 2.0 | 3.0 |

Convolutional Network Components

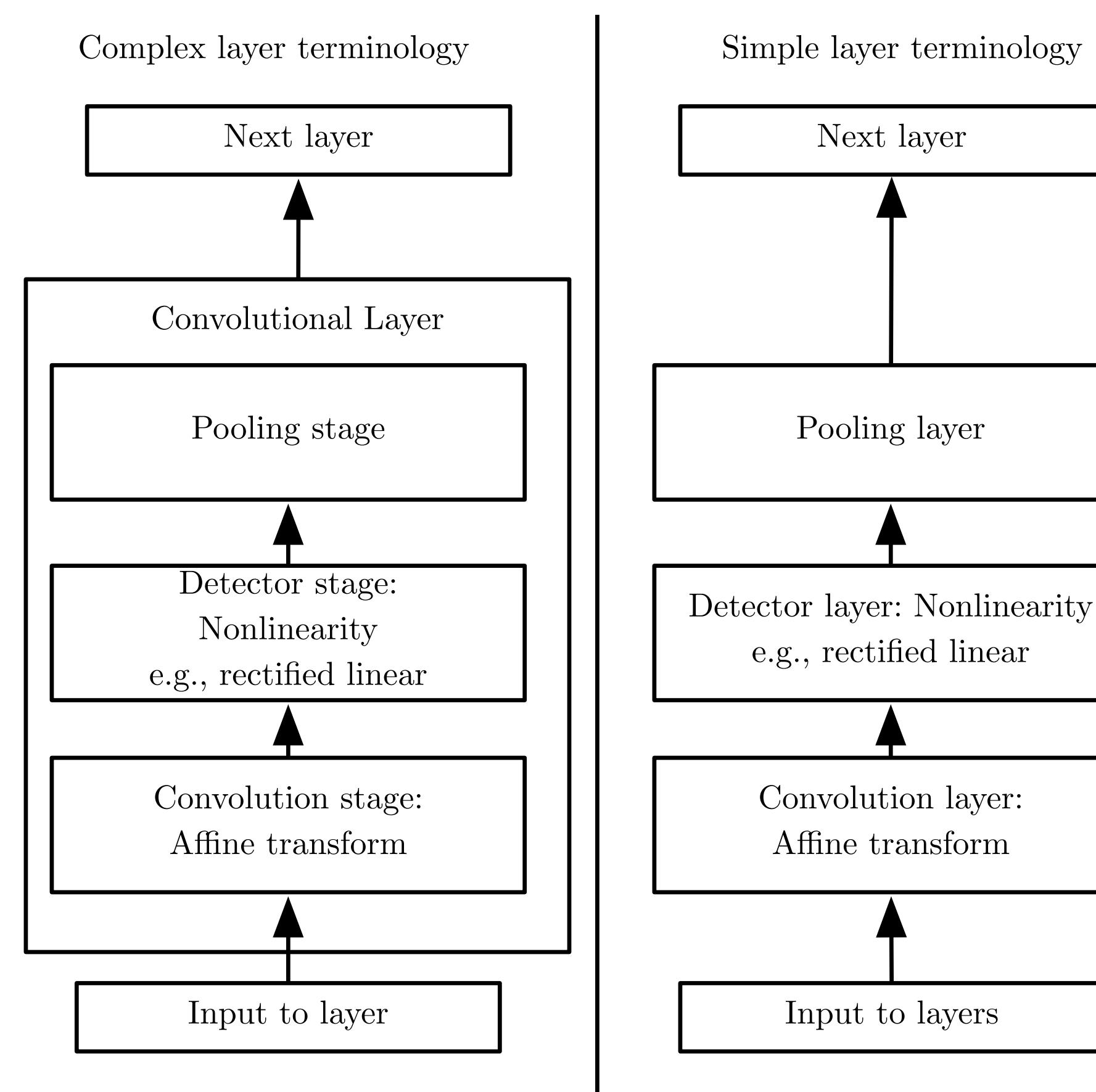


Figure 9.7

(Goodfellow 2016)

Example Classification Architectures

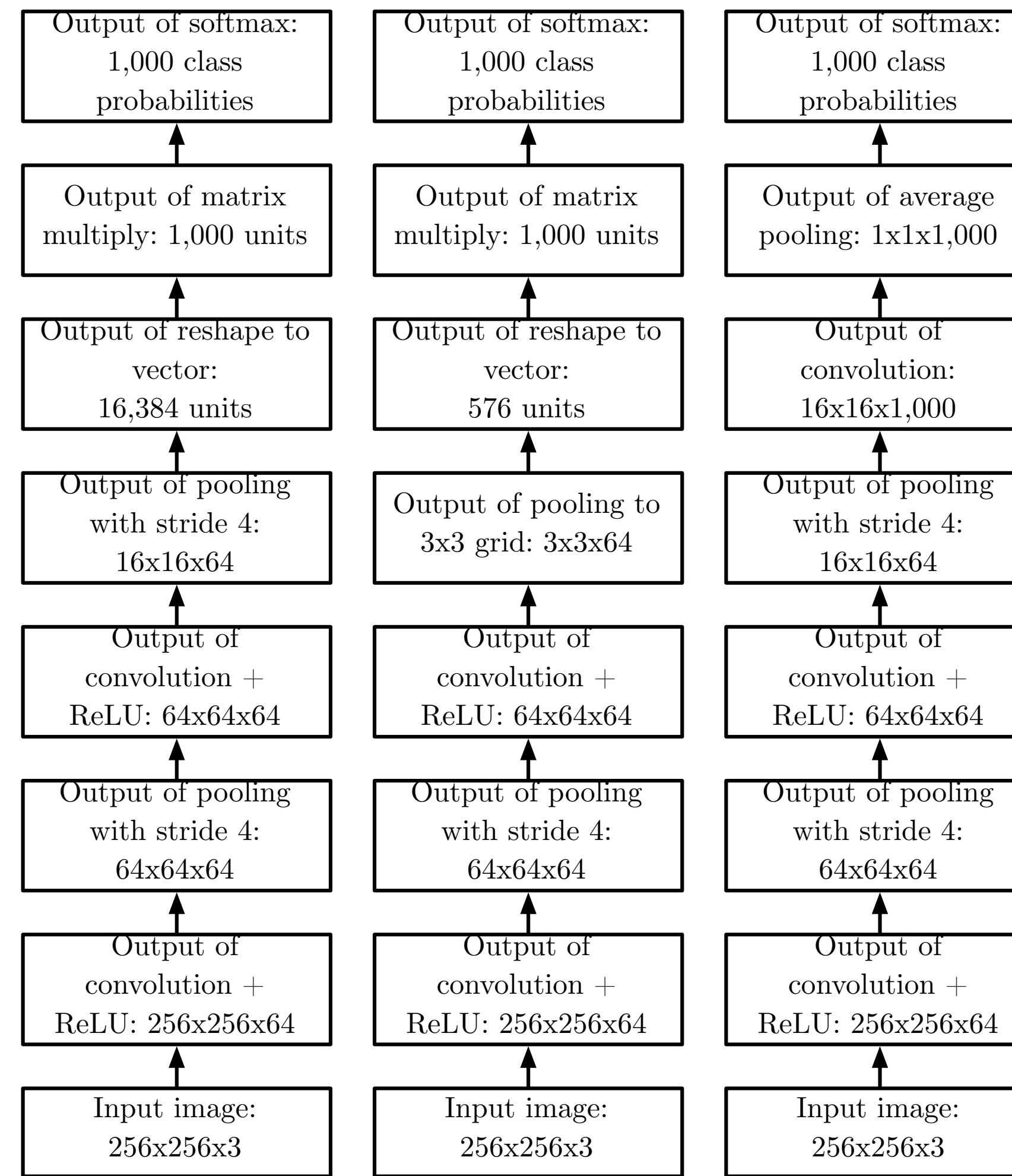


Figure 9.11

It works!

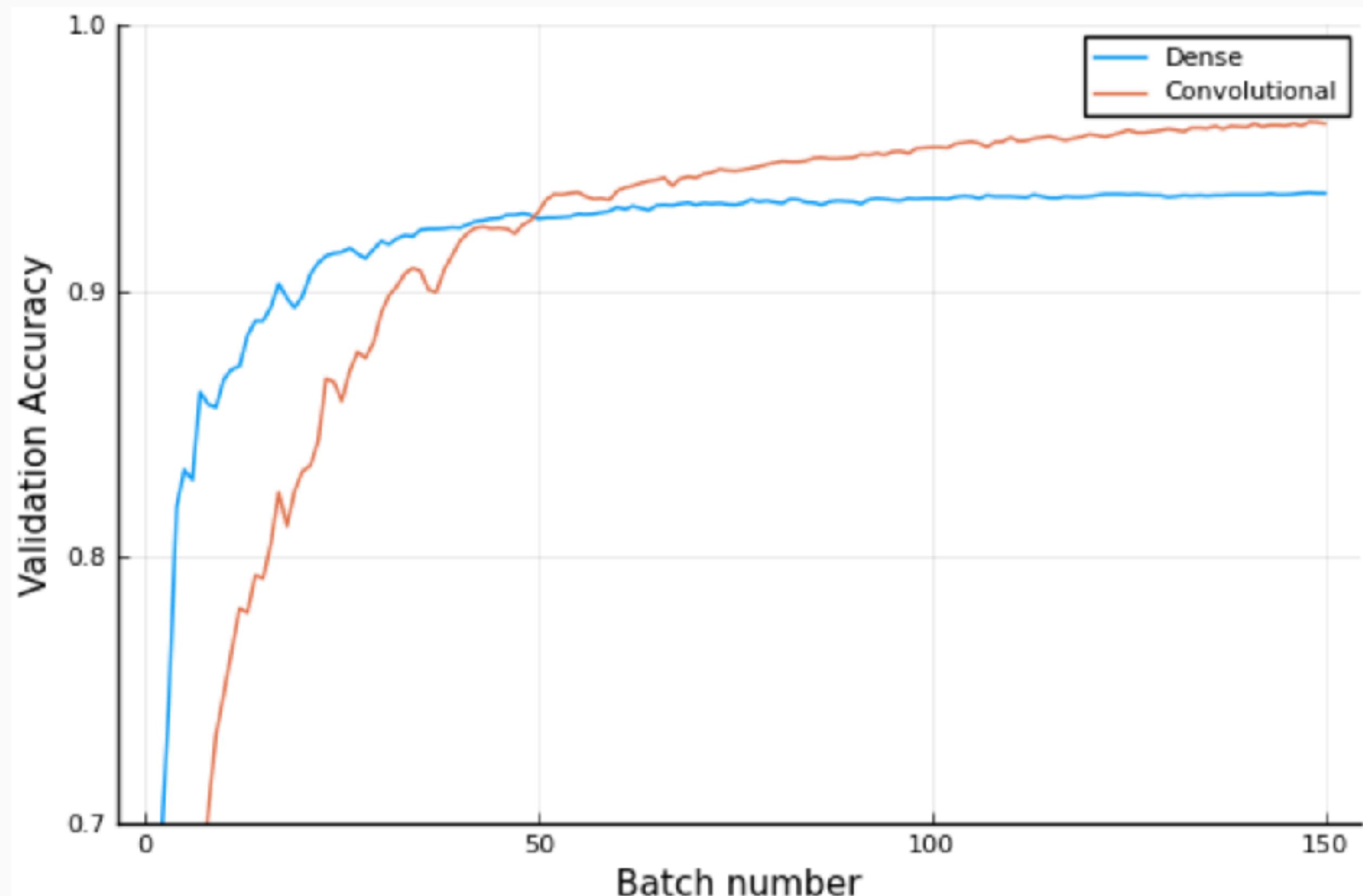
From *MEML*, Ch 5

```
#Julia code
using Flux

model1 = Chain(
    flatten,
    Dense(784, 200, relu),
    Dense(200, 100, tanh),
    Dense(100, 10),
    softmax)

model2 = Chain(
    Conv((5, 5), 1=>8, relu, stride=(1,1), pad=(0,0)),
    MaxPool((2,2)),
    Conv((3, 3), 8=>16, relu, stride=(1,1), pad=(0,0)),
    MaxPool((2,2)),
    flatten,
    Dense(400, 10),
    softmax)
```

Here is a plot of the validation accuracy for both `model1` and `model2` when trained on 5,000 MNIST images (and validated on an additional 5,000). This is with the ADAM optimizer with a learning rate of 0.005 and batch sizes of 1,000.



Mathematics of AI

Emerging themes at end of Module 2

- “Find sparsity”: sparse data representations help generalise
- It’s optimisation all the way down
- Regularisation <-> optimisation <-> sparsity (e.g., compressive sensing)
- Convolutions: the “most useful trick in applied probability”
Fourier transforms ——> CNNs