

Convolutional Neural Networks (CNN)

Dr. Josif Grabocka

ISMLL, University of Hildesheim

Convolutional Neural Networks

The Convolution Operation

- ▶ Generally speaking, convolution is an operation on two functions:

$$s(t) = \int x(a) w(t - a) da$$

- ▶ Often denoted with an asterisk:

$$s(t) = (x * w)(t)$$

- ▶ Example:

- ▶ $x(t)$: a noisy measure the position of a spaceship
- ▶ $w(a)$: relevance of a measurement with age a (Note: $\int w(a) da = 1$)
- ▶ Given a sequence of noisy measurements $x(t), x(t - 1), \dots, x(t - \infty)$, what is the relevance-corrected position $s(t)$?

Convolutions in Deep Learning

- ▶ Terminology: $x(t)$: **Input**, $w(a)$: **Kernel/Filter**, $s(t)$: **Feature Map**
- ▶ Assuming two-dimensional images I and kernels K :

$$S(i, j) = (I * K)(i, j) = \sum_n \sum_m I(i + m, j + n) K(m, n)$$

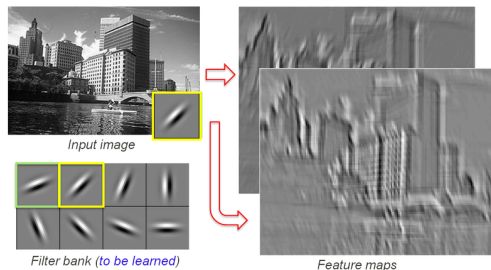
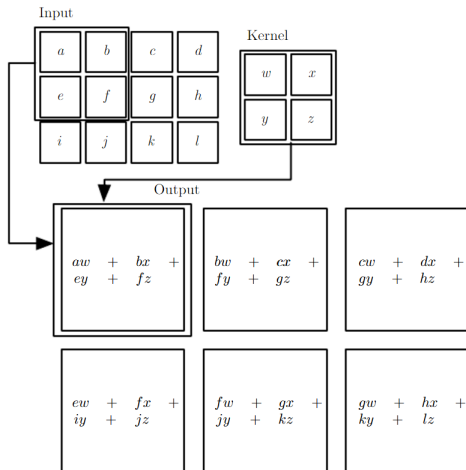


Figure 1: I top left, K bottom-left (green, yellow), S right; Credits: S. Lazebnik

2-D convolution (Source: Goodfellow et al., 2016)



Note: Kernel is **shared** and applied only in **valid** regions

Sparse Interactions/Connectivity

- ▶ Small kernels detect meaningful features (e.g. edges)
- ▶ Reduces memory footprint and computations
- ▶ m inputs, n outputs, kernel of size k reduce the activation complexity from $\mathcal{O}(n \times m)$ to $\mathcal{O}(n \times k)$, for $k \ll m$

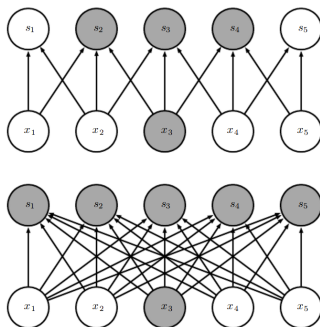


Figure 2: Sparse connectivity; the influence of x_3 ? (Source: Goodfellow et al., 2016)

Multi-layered Sparse Connectivity

- ▶ Stacked convolutions interact with larger portions of the input
- ▶ Capture interactions through sparse connections

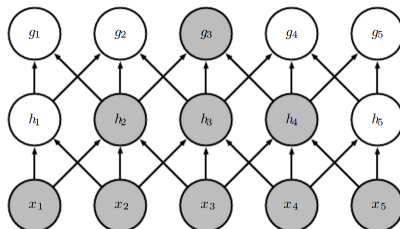


Figure 3: Stacked convolution (Source: Goodfellow et al., 2016)

- ▶ Illustration above:
 - ▶ 2-layer convolution: 26 ops. and 6 weights (assume per-layer kernel)
 - ▶ 2-layer full-nn: 50 ops. and 50 weights

Parameter Sharing or Tied-Weights

- ▶ In a convolutional setup weights are **tied/shared**:

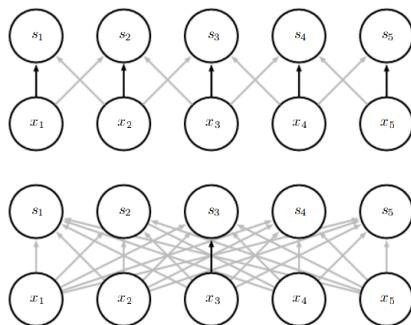


Figure 4: Black arrow represents the same weight (Source: Goodfellow et al., 2016)

- ▶ Achieves a translation-invariant capturing of patterns

Stacked Convolutions with Nonlinear Activations

- ▶ Input image: $V^{(0)} \in \mathbb{R}^{I^{(0)} \times X \times Y}$
- ▶ Kernels: $K^{(1)} \in \mathbb{R}^{I^{(1)} \times C \times M^{(1)} \times N^{(1)}}$, \dots , $K^{(\ell)} \in \mathbb{R}^{I^{(\ell)} \times I^{(\ell-1)} \times M^{(\ell)} \times N^{(\ell)}}$
- ▶ Feature Maps: $V^{(1)} \in \mathbb{R}^{I^{(1)} \times X \times Y}$, \dots , $V^{(\ell)} \in \mathbb{R}^{I^{(\ell)} \times X \times Y}$

$$Z_{i,x,y}^{(1)} = \sum_{c=1}^{I^{(0)}} \sum_{m=1}^{M^{(1)}} \sum_{n=1}^{N^{(1)}} V_{c,x+m-1,y+n-1}^{(0)} K_{i,c,m,n}^{(1)}$$

$$V_{i,x,y}^{(1)} = f\left(Z_{i,x,y}^{(1)}\right), \quad \text{e.g. } f(x) = \max(0, x)$$

$$\vdots$$

$$Z_{i,x,y}^{(\ell)} = \sum_{c=1}^{I^{(\ell-1)}} \sum_{m=1}^{M^{(\ell)}} \sum_{n=1}^{N^{(\ell)}} V_{c,x+m-1,y+n-1}^{(\ell-1)} K_{i,c,m,n}^{(\ell)}$$

$$V_{i,x,y}^{(\ell)} = f\left(Z_{i,x,y}^{(\ell)}\right)$$

Nonlinear Activation of Feature Maps

$$V_{i,x,y}^{(\ell)} = \max \left(0, \sum_{c=1}^{I^{(\ell-1)}} \sum_{m=1}^{M^{(\ell)}} \sum_{n=1}^{N^{(\ell)}} V_{c,x+m-1,y+n-1}^{(\ell-1)} K_{i,c,m,n}^{(\ell)} \right)$$

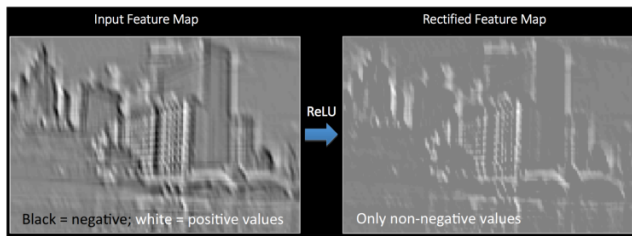


Figure 5: Rectified Feature Map, Credits: Rob Fergus

Strided Convolutions

- Convolve with every s -th position in each dimension:

$$z_{i,x,y}^{(\ell)} = \sum_{c=1}^{I^{(\ell-1)}} \sum_{m=1}^{M^{(\ell)}} \sum_{n=1}^{N^{(\ell)}} V_{c,(x-1)s+m,(y-1)s+n}^{(\ell-1)} K_{i,c,m,n}^{(\ell)}$$

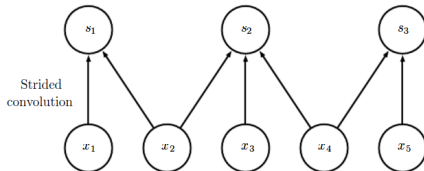


Figure 6: Strided convolutions $s = 2$, (Source: Goodfellow et al., 2016)

Zero Padding - Avoid Size Shrinking

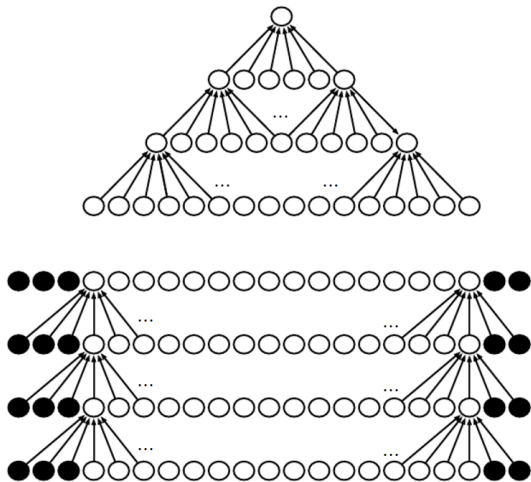


Figure 7: Top: No padding, Bottom: padding (Source: Goodfellow et al., 2016)

Pooling

- ▶ A convolutional network has three stages:
 1. **Convolutions** (in parallel) for multiple kernels
 2. **Nonlinear activations** of the convolutions (ReLU)
 3. **Pooling** (summary statistics)
- ▶ Reduces the dimensionality of the latent representation
- ▶ Ensure invariance to small translations of the input

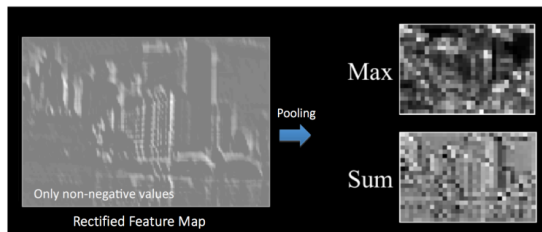


Figure 8: Max and Avg Pooling, Credits: Rob Fergus

Pooling - Translation Invariance Illustration

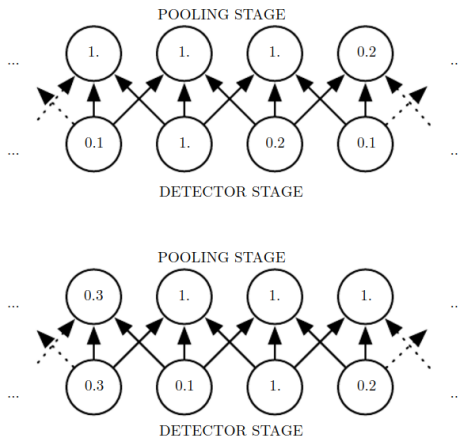


Figure 9: Shifting the input one pixel to the right has smaller effect on the pooling layer, compared to the detector layer (conv+nonlinearity). (Source: Goodfellow et al., 2016)

Pooling - Down-sampling and Strides

Pooling (**max** or **avg**) squared regions of size $\phi \times \phi$ with a stride s :

$$V_{i,x,y}^{(\ell, Pooled)} := \text{Pooling}^{(\ell)} V_{i,\tilde{x},\tilde{y}}^{(\ell)}$$

$$\tilde{x} \in \{(x-1)s+1, \dots, (x-1)s+\phi\}$$

$$\tilde{y} \in \{(y-1)s+1, \dots, (y-1)s+\phi\}$$

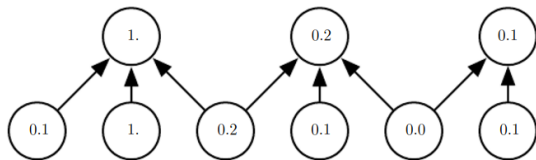
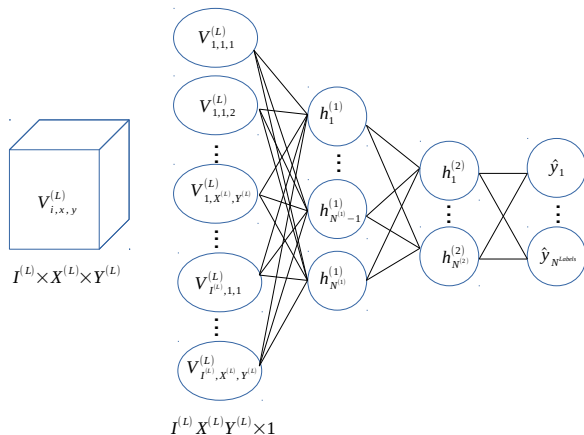


Figure 10: Max-pooling with $s = 2$, $\phi = 3$ (Source: Goodfellow et al., 2016)

Note: For simplicity, we assume $s = \phi$ in the following slides!

Reshaping and Fully Connected Layers



Remember
$$\frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,x,y}^{(L)}} = \sum_i \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial h_i^{(1)}} \frac{\partial h_i^{(1)}}{\partial V_{i,x,y}^{(L)}}$$

CNN - Forward Prediction

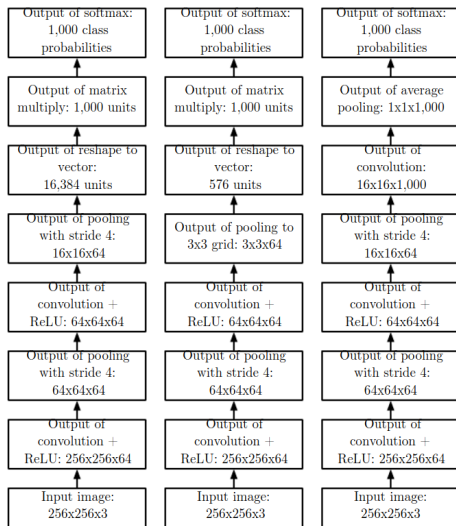
Algorithm 1: Convolutional Neural Network

- 1: {Convolutional steps: $L^{\text{Conv}} \times \{\text{Convolution, Nonlinear, Pooling}\}$ }
- 2: **for** $\ell = 1, \dots, L^{\text{Conv}}$ **do**
- 3:
$$Z_{i,x,y}^{(\ell)} := \sum_{c=1}^{I^{(\ell-1)}} \sum_{m=1}^{M^{(\ell)}} \sum_{n=1}^{N^{(\ell)}} V_{c,(x-1)s^{(\ell,\text{Conv})}+m,(y-1)s^{(\ell,\text{Conv})}+n}^{((\ell-1),\text{Pool})} K_{i,c,m,n}^{(\ell)}$$
- 4:
$$V_{i,x,y}^{(\ell)} := f^{(\ell)}(Z_{i,x,y}^{(\ell)})$$
- 5:
$$V_{i,x,y}^{(\ell,\text{Pool})} := \text{Pooling}^{(\ell)} \left(V_{i,\tilde{x},\tilde{y}}^{(\ell)} \right)$$

$$\tilde{x} \in \{(x-1)s^{(\ell,\text{Pool})}+1, \dots, (x-1)s^{(\ell,\text{Pool})}+\phi^{(\ell)}\}$$

$$\tilde{y} \in \{(y-1)s^{(\ell,\text{Pool})}+1, \dots, (y-1)s^{(\ell,\text{Pool})}+\phi^{(\ell)}\}$$
- 6: {Fully connected layers}
- 7:
$$h^{(0)} := \left[V_{1,1,1}^{(L^{\text{Conv}})}, \dots, V_{I^{(L)}, X^{(L)}, Y^{(L)}}^{(L^{\text{Conv}})} \right]$$
- 8: **for** $\ell = 1, \dots, L^{\text{Full}}$ **do**
- 9:
$$h_{\cdot}^{(\ell)} = f^{(\ell)}(W_{\cdot,\cdot}^{(\ell)} h_{\cdot}^{(\ell-1)} + b_{\cdot}^{(\ell)})$$
- 10: **return** $\hat{y} := h_{\cdot}^{L^{\text{Full}}}$

Example Architectures (Source: Goodfellow et al., 2016)



Gradients

- ▶ Convolutions and pooling are computational graph nodes
- ▶ Apply the standard back-propagation for computational graphs
- ▶ Remember a convolution:

$$Z_{i,x,y}^{(\ell)} = \sum_{c=1}^{I^{(\ell-1)}} \sum_{m=1}^{M^{(\ell)}} \sum_{n=1}^{N^{(\ell)}} V_{c,(x-1)s+m,(y-1)s+n}^{(\ell-1)} K_{i,c,m,n}^{(\ell)}$$

- ▶ Given $\frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,x,y}^{(\ell)}}$, define $G_{i,x,y}^{(\ell)} := \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial Z_{i,x,y}^{(\ell)}} = \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,x,y}^{(\ell)}} \left(f^{(\ell)'}(Z_{i,x,y}^{(\ell)}) \right)$
- ▶ Yielding:

$$\frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial K_{i,c,m,n}^{(\ell)}} = \sum_{x=1}^{X^{(\ell)}} \sum_{y=1}^{Y^{(\ell)}} G_{i,x,y}^{(\ell)} V_{c,(x-1)s+m,(y-1)s+n}^{(\ell-1)}$$

Gradients (II)

- ▶ We need the gradient w.r.t. $V_{:, :, :}^{(\ell-1)}$ to propagate the error down:

$$\frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i, x, y}^{(\ell-1)}} = \sum_{(x'-1)s+p=x}^{x', p} \sum_{(y'-1)s+q=y}^{y', q} \sum_{c=1}^{I^{(\ell)}} K_{c, i, p, q}^{(\ell)} G_{c, x', y'}^{(\ell)}$$

- ▶ Gradients of pooling are simpler, e.g. for max layer and $s = \phi$:

$$\frac{\partial V_{i, x, y}^{(\ell, Pooled)}}{\partial V_{i, \tilde{x}, \tilde{y}}^{(\ell)}} = \begin{cases} 1 & \text{if } (\tilde{x}, \tilde{y}) = \underset{\substack{\tilde{x}^* \in \{(x-1)s+1, \dots, (x-1)s+\phi\} \\ \tilde{y}^* \in \{(y-1)s+1, \dots, (y-1)s+\phi\}}}]{\text{argmax}} V_{i, \tilde{x}^*, \tilde{y}^*}^{(\ell)} \\ 0 & \text{else} \end{cases}$$

CNN - Back-Propagation

Algorithm 2: CNN's parameters gradients (Only Convolutional Layers)

$$1: \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,x,y}^{(L^{\text{Conv}}, \text{Pool})}} = \sum_i \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial h_i^{(1)}} \frac{\partial h_i^{(1)}}{\partial V_{i,x,y}^{(L^{\text{Conv}}, \text{Pool})}}$$

2: **for** $\ell = L^{\text{Conv}}, \dots, 1$ **do**

$$3: \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,\tilde{x},\tilde{y}}^{(\ell)}} = \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,x,y}^{(\ell, \text{Pool})}} \times \begin{cases} 1 & (\tilde{x}, \tilde{y}) = \underset{\substack{\tilde{x}^* \in \{(x-1)s+1, \dots, (x-1)s+\phi\} \\ \tilde{y}^* \in \{(y-1)s+1, \dots, (y-1)s+\phi\}}}]{\text{argmax}} V_{i,\tilde{x}^*,\tilde{y}^*}^{(\ell)} \\ 0 & \text{else} \end{cases}$$

$$4: G_{i,\tilde{x},\tilde{y}}^{(\ell)} := \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,\tilde{x},\tilde{y}}^{(\ell)}} \left(f^{(\ell)'}(Z_{i,\tilde{x},\tilde{y}}^{(\ell)}) \right)$$

$$5: \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial K_{i,c,m,n}^{(\ell)}} = \sum_{\tilde{x}=1}^{X^{(\ell)}} \sum_{\tilde{y}=1}^{Y^{(\ell)}} G_{i,\tilde{x},\tilde{y}}^{(\ell)} V_{c,(\tilde{x}-1)s+m,(\tilde{y}-1)s+n}^{(\ell-1, \text{Pool})}$$

$$6: \frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial V_{i,\tilde{x},\tilde{y}}^{(\ell-1, \text{Pool})}} = \sum_{\substack{x',p \\ (x'-1)s+p=\tilde{x}}} \sum_{\substack{y',q \\ (y'-1)s+q=\tilde{y}}} \sum_{c=1}^{I^{(\ell)}} K_{c,i,p,q}^{(\ell)} G_{c,x',y'}^{(\ell)}$$

7: **return** $\frac{\partial \mathcal{L}(\mathcal{Y}, \hat{\mathcal{Y}})}{\partial K}$