# Computing Convolutions
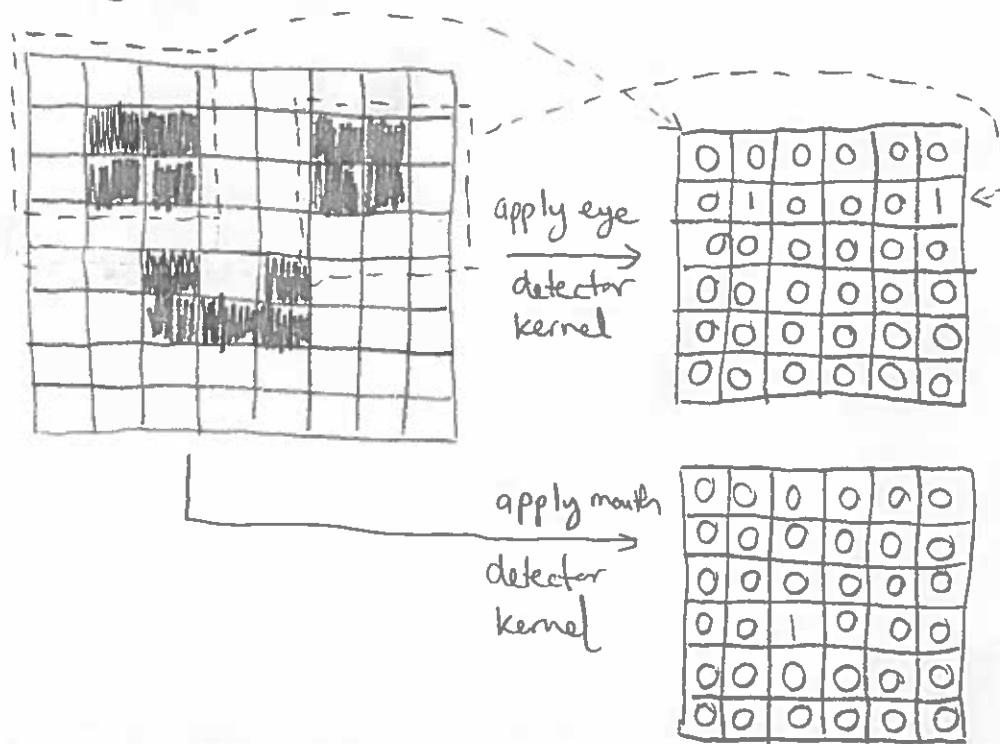
① When we run the happy face detector, we need to apply each kernel to each 3×3 window of the image:



apply eye detector kernel →

apply mouth detector kernel →

② By "apply a kernel to a window," we mean that we compute:

Hadamard (elementwise) product



$$\text{sum} \left( \boxed{\text{image window}} \ \oplus \ \boxed{\text{convolution kernel}} \right)$$

image window          convolution kernel

# COMPUTING CONVOLUTIONS

③ How can we do this efficiently (i.e. leverage matrix multiplication libraries as much as possible)?

Let's consider a simpler example, in which we have a $3 \times 3$ "image":

$$\begin{array}{|c|c|c|} \hline 1 & 1 & 0 \\ \hline 0 & 0 & 1 \\ \hline 1 & 0 & 1 \\ \hline \end{array}$$

and a single $2 \times 2$ kernel:

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 0 & -1 \\ \hline \end{array}$$

④ To "convolve" the image with the kernel (with no padding and stride 1), we want to compute the matrix:

$$\left[ \begin{array}{c|c}
\text{sum}\left( \begin{array}{|c|c|}\hline 1 & 1 \\\hline 0 & 0 \\\hline\end{array} \oplus \begin{array}{|c|c|}\hline 1 & 2 \\\hline 0 & -1 \\\hline\end{array} \right) & \text{sum}\left( \begin{array}{|c|c|}\hline 1 & 0 \\\hline 0 & 1 \\\hline\end{array} \oplus \begin{array}{|c|c|}\hline 1 & 2 \\\hline 0 & -1 \\\hline\end{array} \right) \\
\hline
\text{sum}\left( \begin{array}{|c|c|}\hline 0 & 0 \\\hline 1 & 0 \\\hline\end{array} \oplus \begin{array}{|c|c|}\hline 1 & 2 \\\hline 0 & -1 \\\hline\end{array} \right) & \text{sum}\left( \begin{array}{|c|c|}\hline 0 & 1 \\\hline 0 & 1 \\\hline\end{array} \oplus \begin{array}{|c|c|}\hline 1 & 2 \\\hline 0 & -1 \\\hline\end{array} \right)
\end{array} \right]$$

$$= \begin{array}{|c|c|} \hline 3 & 0 \\ \hline 0 & 1 \\ \hline \end{array}$$

# Computing Convolutions

⑤ We can express these computations as a matrix multiplication by expressing the kernel as a row vector:

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 0 & -1 \\ \hline \end{array} \longrightarrow [1 \ 2 \ 0 \ -1]$$

and the image windows as columns of a matrix:

$$\begin{array}{|c|c|c|} \hline 1 & 1 & 0 \\ \hline 0 & 0 & 1 \\ \hline 1 & 0 & 1 \\ \hline \end{array} \longrightarrow \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

⑥ When we multiply the row vector by the column matrix, we get:

$$[1 \ 2 \ 0 \ -1] \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix} = [3 \ 0 \ 0 \ 1]$$

$$sum \left( \begin{array}{|c|c|} \hline 1 & 1 \\ \hline 0 & 0 \\ \hline \end{array} \oplus \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 0 & -1 \\ \hline \end{array} \right)$$

which we can reshape into a 2x2 matrix to get the desired result:

$$[3 \ 0 \ 0 \ 1] \longrightarrow \begin{array}{|c|c|} \hline 3 & 0 \\ \hline 0 & 1 \\ \hline \end{array}$$

# Computing Convolutions

⑦ It's not too hard to do this with multiple kernels at once (e.g. an eye detector and a mouth detector).

Let's add a kernel to our example, so we have now a 3×3 image:

$$\begin{array}{|c|c|c|} \hline 1 & 1 & 0 \\ \hline 0 & 0 & 1 \\ \hline 1 & 0 & 1 \\ \hline \end{array}$$

and two 2×2 kernels:

$$\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 0 & -1 \\ \hline \end{array} \quad \text{and} \quad \begin{array}{|c|c|} \hline 3 & 0 \\ \hline -2 & 4 \\ \hline \end{array}$$

---

⑧ Convolving the image with these kernels (with no padding and stride 1), we get:



$$= \left[ \begin{array}{|c|c|} \hline 3 & 0 \\ \hline 0 & 1 \\ \hline \end{array} \; , \; \begin{array}{|c|c|} \hline 3 & 7 \\ \hline -2 & -2 \\ \hline \end{array} \right]$$

# COMPUTING CONVOLUTIONS

⑨ We can perform the same trick, expressing each kernel as a row in a "kernel-row" matrix, and each image window as a column in a "window column" matrix. Then we multiply the matrices to obtain our desired quantities:

$$
\begin{array}{c}
\text{kernel 1} \rightarrow \\
\text{kernel 2} \rightarrow
\end{array}
\underbrace{\begin{bmatrix} 1 & 2 & 0 & -1 \\ 3 & 0 & -2 & 4 \end{bmatrix}}_{\text{kernel-row matrix}}
\overbrace{\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{bmatrix}}^{\text{window-column matrix}}
=
\begin{bmatrix} 3 & 0 & 0 & 1 \\ 3 & 7 & -2 & -2 \end{bmatrix}
$$

(↑ image window 1, ↑ image window 4)

which we can reshape into a 2×2×2 tensor to get the desired result:

$$
\begin{bmatrix} 3 & 0 & 0 & 1 \\ 3 & 7 & -2 & -2 \end{bmatrix}
\rightarrow
\left[ \begin{array}{cc} 3 & 0 \\ 0 & 1 \end{array} \, , \, \begin{array}{cc} 3 & 7 \\ -2 & -2 \end{array} \right]
$$

# COMPUTING CONVOLUTIONS

⑩ There are a couple additional complications. While grayscale images are stored as matrices, color images are often represented by order-3 tensors — one matrix for each red, green, and blue "channel", e.g.

$$\begin{bmatrix} \begin{array}{|c|c|c|} \hline 1 & 0 & 1 \\ \hline 0 & 1 & 0 \\ \hline 1 & 0 & 1 \\ \hline \end{array} \underset{\text{red channel}}{} , \begin{array}{|c|c|c|} \hline 1 & 1 & 0 \\ \hline 0 & 0 & 1 \\ \hline 1 & 0 & 1 \\ \hline \end{array} \underset{\text{green channel}}{} , \begin{array}{|c|c|c|} \hline 0 & 0 & 1 \\ \hline 0 & 1 & 0 \\ \hline 1 & 0 & 0 \\ \hline \end{array} \underset{\text{blue channel}}{} \end{bmatrix}$$

⑪ In this case, each kernel is also an order-3 tensor, one matrix for each channel, e.g.:

$$\begin{bmatrix} \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 0 & 0 \\ \hline \end{array} \underset{\substack{\text{red} \\ \text{kernel}}}{} , \begin{array}{|c|c|} \hline 0 & 0 \\ \hline 1 & 0 \\ \hline \end{array} \underset{\substack{\text{green} \\ \text{kernel}}}{} , \begin{array}{|c|c|} \hline -1 & 0 \\ \hline 0 & 0 \\ \hline \end{array} \underset{\substack{\text{blue} \\ \text{kernel}}}{} \end{bmatrix}$$

⑫ The convolution is:

| | |
|---|---|
| $\text{Sum}\left(\boxplus \oplus \boxplus\right) + \text{Sum}\left(\boxplus \oplus \boxplus\right)$ $+ \text{Sum}\left(\boxplus \oplus \boxplus\right)$ | $\text{Sum}\left(\boxplus \oplus \boxplus\right) + \text{Sum}\left(\boxplus \oplus \boxplus\right)$ $+ \text{Sum}\left(\boxplus \oplus \boxplus\right)$ |
| | |

$$= \begin{array}{|c|c|} \hline 1 & 0 \\ \hline 1 & 0 \\ \hline \end{array}$$

⑬ Ok, no problem! We can still compute this as the multiplication of a kernel-row matrix with a window-column matrix: we just need to concatenate the channels.



$$\text{kernel 1} \rightarrow \begin{bmatrix} \underset{\text{red}}{\begin{array}{|c|c|}\hline 1 & 0 \\\hline 0 & 0 \\\hline\end{array}} , & \underset{\text{green}}{\begin{array}{|c|c|}\hline 0 & 0 \\\hline 1 & 0 \\\hline\end{array}} , & \underset{\text{blue}}{\begin{array}{|c|c|}\hline -1 & 0 \\\hline 0 & 0 \\\hline\end{array}} \end{bmatrix}$$

$$\overset{\text{kernel 1}}{=} \begin{bmatrix} \underbrace{1\ 0\ 0\ 0}_{\text{red}}\ \underbrace{0\ 0\ 1\ 0}_{\text{green}}\ \underbrace{-1\ 0\ 0\ 0}_{\text{blue}} \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & 1 & 0 \end{bmatrix}$$
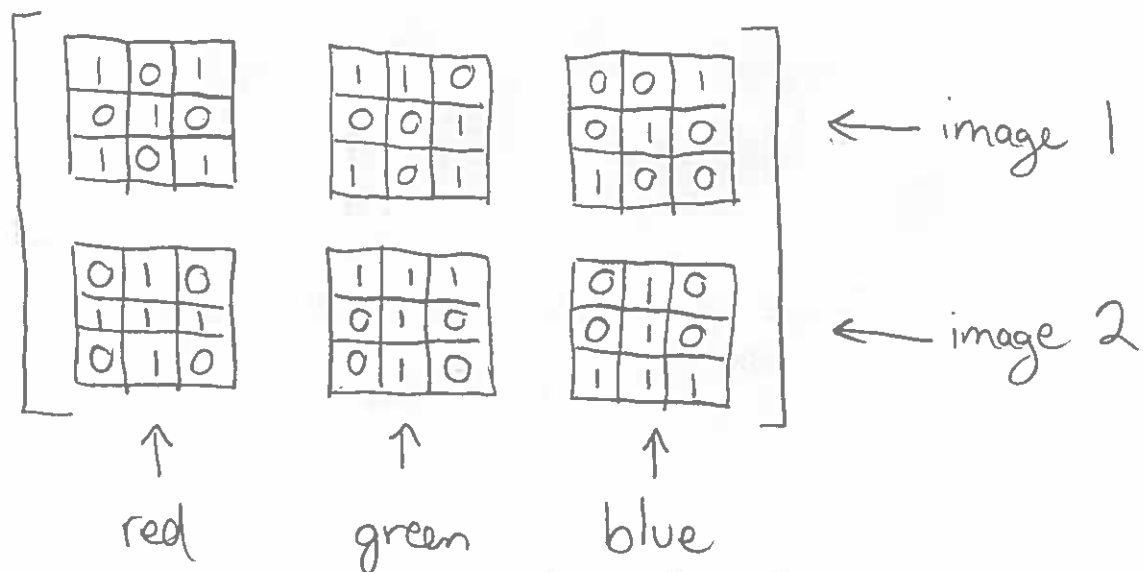
⑭ The resulting matrix can then be reshaped into a $1 \times 2 \times 2$ tensor to get the result of convolution:

↑ num kernels

↑ resulting image size

$$[\; 1 \quad 0 \quad 1 \quad 0\;] \rightarrow \left[\begin{array}{|c|c|}\hline 1 & 0 \\\hline 1 & 0 \\\hline\end{array}\right]$$

⑮ All right. That's almost, but not quite, everything. The last complication is, in minibatch training, we're usually not convolving one image at a time, but actually a batch of images.

So the input to our convolution might actually be an order-4 tensor like this:



← image 1

← image 2

↑ red      ↑ green      ↑ blue

This tensor has shape $\underbrace{2}_{\text{num images}} \times 3 \times \overbrace{3 \times 3}^{\text{image size}}$.

↳ RGB

(16) This is a straightforward (ish) extension. The kernels haven't changed, so the kernel-row matrix doesn't change. We just need to add the other image's windows as additional columns of the window-column matrix.

kernel 1 → [ red green blue | image 1 (win1 win2 win3 win4) | image 2 (win1 win2 win3 win4) ] (red, green, blue rows)

$$= \begin{bmatrix} 1 & 0 & 0 & 1 & 0 & 0 & 1 & 0 & -1 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} \text{image 1 (win1 win2 win3 win4)} & \text{image 2 (win1 win2 win3 win4)} \\ \vdots & \vdots \end{bmatrix}$$

The window-column matrix (image 1: win1 win2 win3 win4, image 2: win1 win2 win3 win4) with rows grouped red, green, blue:

red:
| | win1 | win2 | win3 | win4 | win1 | win2 | win3 | win4 |
|---|---|---|---|---|---|---|---|---|
| | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 |
| | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 |
| | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 |
| | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |

green:
| | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 |
| | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| | 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |

blue:
| | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 |
| | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |
| | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 |
| | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 |

$$= \begin{bmatrix} 1 & 0 & 1 & 0 & 1 & 2 & 2 & 1 \end{bmatrix}$$

which can be reshaped into a $2 \times 1 \times 2 \times 2$ tensor (num images, num kernels, image size) to get the result of convolution:

image 1:
| 1 | 0 |
|---|---|
| 1 | 0 |

image 2:
| 1 | 2 |
|---|---|
| 2 | 1 |