

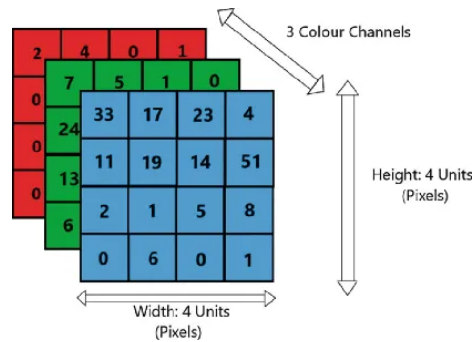
EMAT31530: Convolutional neural networks

Laurence Aitchison

So far, we've seen lots of material on how to train NNs using backprop, and how to use PyTorch. But to actually apply NNs to e.g. image datasets (for instance for classifying the objects in the image), we need to consider architectures specifically adapted to images. Perhaps the key operation is “convolution”, which is the basis for all the classic and many of the modern networks for images. Networks that involve convolutional layers are called “convolutional neural networks”, and that's what we'll look at this week!

1 Images and feature maps as tensors

To understand convolutions, the first step is to understand how we represent images in deep learning. We represent images as big 3D tensors, with shape $3 \times \text{height} \times \text{width}$,



The three colour channels represent the red, green and blue channels in the image (i.e. the amount of red, green or blue at that pixel). Alternatively, a greyscale image is a $1 \times \text{height} \times \text{width}$ image.

A feature map is very similar to an image. The key difference is that a feature map has much larger number of channels (usually around $100 - 500$), rather than 3 in a colour image or 1 in a greyscale image. So a feature map is of shape $\text{channels} \times \text{height} \times \text{width}$, which is usually abbreviated to $C \times H \times W$. You can imagine that in a feature map, we have assigned a feature vector to every pixel in every image. That contrasts to fully-connected networks, where we'd just have one feature vector associated with every image.

In practice, we usually work with a batch of N images. Thus, we end up working with tensors of size minibatch size \times channels \times height \times width, which is usually abbreviated to $N \times C \times H \times W$.

Convolutional layers take images/feature maps as inputs, and returns images/feature maps as outputs. We'll see how they operate in the next section.

2 Convolutions

2.1 Intuition for convolutions in the simplest possible case

To start, we draw a diagram for the simplest possible deep-learning style convolution. Specifically, consider a convolution with one input channel and one output channel, and with weights,

0	1	2
2	2	0
0	1	2

We slide these weights over the input image/feature map. At each location, we calculate the product between the weight and the image pixel. We sum all these products to get the output in the current location.

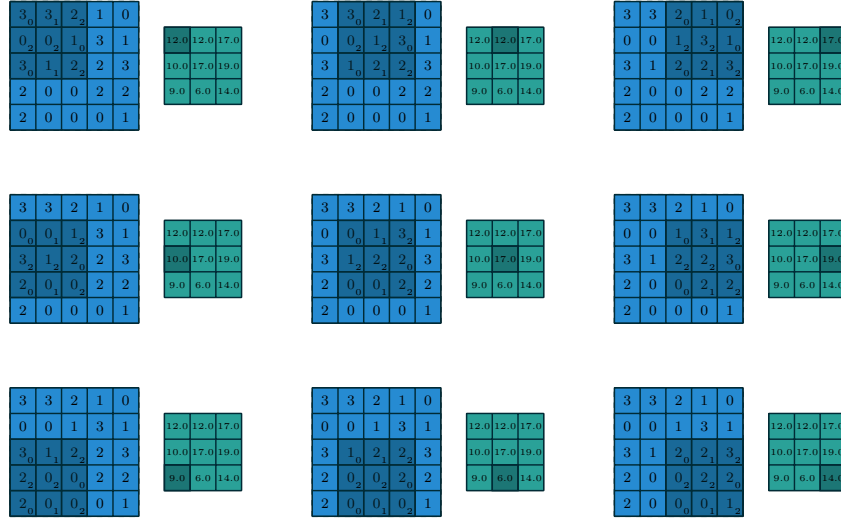


Figure 1: Computing the output of a simple convolutional layer. The blue 5×5 square shows the inputs, the green 3×3 square shows the output. The darker 3×3 square in the input shows the current location of the patch of weights as it slides over the image. This diagram along with others in this pdf are taken from “A guide to convolution arithmetic for deep learning”, by Vincent Dumoulin and Francesco Visin. This is a great article to check out if you’re interested in filling in any gaps.

Specifically, the formula for this simple convolution is,

$$a_{x,y} = \sum_{\delta_x=-1}^1 \sum_{\delta_y=-1}^1 h_{x+\delta_x, y+\delta_y} W_{\delta_x, \delta_y} \quad (1)$$

here:

- $a_{x,y}$ is the output image / feature map. In this example, the outputs are size 3×3 . So the indices x and y represent location in the output image, and range from 1 to 3.
- $h_{x+\delta_x, y+\delta_y}$ is the input image / feature map. In this example, the inputs are size 5×5 . So the indices $x + \delta_x$ and $y + \delta_y$ represent location on the input image, and range from 0 to 4.
- W_{δ_x, δ_y} is the convolutional weights (in grey in the first image, or in the 3×3 darker patches in the inputs in Fig. 1). In this example, the weights are size 3×3 . So the indices δ_x and δ_y represent the location within the 3×3 weight matrix, and range from -1 to 1 .

Note that here, when we say that “an index ranges from -1 to 1 ”, we’re being flexible for the purposes of making the maths easy. In Python/PyTorch, indices always range from 0 to $(\text{length} - 1)$.

There are a couple of ways of interpreting a convolution in terms of matrices and vectors. Both of these interpretations are easier to understand if we consider a 1D convolution, where instead of the inputs and outputs being images with 2 spatial dimensions, they just have one spatial dimension, x . In that case, Eq. (1) becomes,

$$a_x = \sum_{\delta=-1}^1 h_{x+\delta} W_{\delta}. \quad (2)$$

Writing out the individual components (for a length-five input and a length five output),

$$a_1 = h_0 W_{-1} + h_1 W_0 + h_2 W_1 \quad (3a)$$

$$a_2 = h_1 W_{-1} + h_2 W_0 + h_3 W_1 \quad (3b)$$

$$a_3 = h_2 W_{-1} + h_3 W_0 + h_4 W_1 \quad (3c)$$

2.2 Interpretation 1: linear operation on patches

Notice that the components in Eq. (2) or Eq. (3) can be written as a big matrix operation,

$$a_x = \begin{pmatrix} h_{x-1} & h_x & h_{x+1} \end{pmatrix} \begin{pmatrix} W_{-1} \\ W_0 \\ W_{+1} \end{pmatrix} = \mathbf{p}_x \mathbf{W} \quad (4)$$

Here, \mathbf{p}_x is the inputs in the patch around h_x concatenated together

$$\mathbf{p}_x = (h_{x-1} \quad h_x \quad h_{x+1}), \quad (5)$$

and \mathbf{W} is the weights for all parts of the patch concatenated together,

$$\mathbf{W} = \begin{pmatrix} W_{-1} \\ W_0 \\ W_1 \end{pmatrix}. \quad (6)$$

2.3 Interpretation 2: linear operations on the input image

There's an alternative interpretation, which doesn't require us to extract patches, \mathbf{p}_x , from the input image. Specifically, we can write out the result of Eq. (3) as an alternative big matrix-vector operation, with a structured weight matrix,

$$(a_1 \quad a_2 \quad a_3) = (h_0 \quad h_1 \quad h_2 \quad h_3 \quad h_4) \begin{pmatrix} W_{-1} & 0 & 0 \\ W_0 & W_{-1} & 0 \\ W_1 & W_0 & W_{-1} \\ 0 & W_1 & W_0 \\ 0 & 0 & W_1 \end{pmatrix} \quad (7)$$

Note that the input and output vectors here are just the standard single-channel input and output vectors, \mathbf{h} and \mathbf{a} , so we can write this expression as,

$$\mathbf{a} = \mathbf{h} \begin{pmatrix} W_{-1} & 0 & 0 \\ W_0 & W_{-1} & 0 \\ W_1 & W_0 & W_{-1} \\ 0 & W_1 & W_0 \\ 0 & 0 & W_1 \end{pmatrix}. \quad (8)$$

In fact, all convolutions can be understood as matrix multiplications with structured weights. That interpretation is okay in this simple setting, with one input and output channel, and one spatial dimension, but becomes quite complex as soon as you consider multiple input/output channels, or more than one spatial dimension.

2.4 Convolutions with multiple input and output channels

In the previous section, we considered the simplest setting for convolutions, with a single input channel, and a single output channel. However, in practice in deep learning, we almost always have multiple input and multiple output channels. What happens in that setting? Well, it is almost exactly the same as in the single-input and single-output channel settings, except:

- Instead of an scalar at each location in the input and output (i.e. $h_{x+\delta x, y+\delta y}$ and $a_{x,y}$), there's a whole feature *vector* at each location, $\mathbf{h}_{x+\delta x, y+\delta y}$ and $\mathbf{a}_{x,y}$.

- Instead of just a scalar weight at each location in the patch (i.e. W_{δ_x, δ_y}) that there's a whole weight matrix, $\mathbf{W}_{\delta_x, \delta_y}$.

Overall, that means a 2D convolution, with multiple input and output channels can be written,

$$\mathbf{a}_{x,y} = \sum_{\delta_x=-1}^1 \sum_{\delta_y=-1}^1 \mathbf{h}_{x+\delta_x, y+\delta_y} \mathbf{W}_{\delta_x, \delta_y}. \quad (9)$$

Here,

- $\mathbf{a}_{x,y}$ is a row-vector of length C_{out} . There is one vector at each spatial location in the output.
- $\mathbf{h}_{x+\delta_x, y+\delta_y}$ is a row-vector of length C_{in} . There is one vector at each spatial location in the input.
- $\mathbf{W}_{\delta_x, \delta_y}$ is a matrix of size $C_{\text{in}} \times C_{\text{out}}$. There is one of these matrices at each location in the patch.

Interpretations 1 + 2 can also be extended to the multi-dimensional case, by directly swapping the scalar inputs/outputs/weights for their vector/matrix counterparts, but we're only going to look at extending interpretation 1.

2.5 Interpretation 1 in the multidimensional setting

A convolution with 1 spatial dimension, but multiple input and output channels looks like,

$$\mathbf{a}_x = \sum_{\delta=-1}^1 \mathbf{h}_{x+\delta} \mathbf{W}_{\delta}. \quad (10)$$

Notice that Eq. (10) can be written as a big matrix operation,

$$\mathbf{a}_x = (\mathbf{h}_{x-1} \quad \mathbf{h}_x \quad \mathbf{h}_{x+1}) \begin{pmatrix} \mathbf{W}_{-1} \\ \mathbf{W}_0 \\ \mathbf{W}_{+1} \end{pmatrix} = \mathbf{p}_x \mathbf{W} \quad (11)$$

Here, \mathbf{p}_x represents inputs in the patch around \mathbf{h}_x concatenated together. Remember that \mathbf{h}_x is a row vector of size C_{in} , so \mathbf{p}_x is a bigger row-vector, of length $(3C_{\text{in}})$.

$$\mathbf{p}_x = (\mathbf{h}_{x-1} \quad \mathbf{h}_x \quad \mathbf{h}_{x+1}) \quad (12)$$

Moreover, \mathbf{W} is a joint matrix for all locations in the patch,

$$\mathbf{W} = \begin{pmatrix} \mathbf{W}_{-1} \\ \mathbf{W}_0 \\ \mathbf{W}_1 \end{pmatrix} \quad (13)$$

Remember that \mathbf{W}_{δ} is of size $C_{\text{in}} \times C_{\text{out}}$. So \mathbf{W} is a matrix of size $(3C_{\text{in}}) \times C_{\text{out}}$. Thus, \mathbf{W} is the right shape to take whole patches, $\mathbf{p}_{x,y}$ (which are vectors of length $3C_{\text{in}}$) as inputs.

2.6 Shapes in PyTorch in practice

Finally, it is worth mentioning the shape of the inputs and weights in actual convolutions in PyTorch.

- The inputs and outputs are standard feature maps, of size $N \times C_{\text{in}} \times H_{\text{in}} \times W_{\text{in}}$ and $N \times C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}$. Note that the inputs and outputs can be of different sizes (in Fig. 1 we had $5 = H_{\text{in}} = W_{\text{in}}$ and $3 = H_{\text{out}} = W_{\text{out}}$). Finally, N is the minibatch size.
- The weights are of size $C_{\text{out}} \times C_{\text{in}} \times H_{\text{kernel}} \times W_{\text{kernel}}$, where H_{kernel} and W_{kernel} are the size of the convolutional patch (or “kernel”); (in Fig. 1 we had $3 = H_{\text{kernel}} = W_{\text{kernel}}$).

2.7 Convolutions in (Engineering) Maths [Non-Examinable]

If you have not encountered convolutions before, you can ignore this section. However, if you have encountered convolutions before in an (Engineering) maths context, you might be slightly confused, as convolutions are typically defined as,

$$a(x) = \int_{-\infty}^{\infty} d\delta' h(x - \delta') w(\delta') \quad (14)$$

Note that here, δ is a variable, not the Kronecker/Dirac delta. This is an operator that takes two functions (h and w) as input, and returns a function, a as output. This doesn't look much like the AI-version of convolution. We can make it look a bit more similar by setting, $\delta = -\delta'$

$$a(x) = \int_{-\infty}^{\infty} d\delta h(x + \delta) w(-\delta). \quad (15)$$

If we write down a discrete version of this expression, we get something that looks a lot like an AI-style convolution,

$$a_x = \sum_{\delta=-\infty}^{\infty} h_{x+\delta} W_{-\delta}. \quad (16)$$

There are two key differences:

- Here, δ ranges from $-\infty$ to ∞ . In contrast, in AI (Eq. 2), δ has a small range, e.g. often from -1 to 1 .
- Here, there's a $-\delta$ in $W_{-\delta}$, while in AI (Eq. 2) there's no minus sign so we just have W_{δ} . It turns out that this doesn't matter in AI, as all the elements of the weight are learned.

2.8 Strides and Padding

When practically using convolutions in deep learning, there are a few options to set. In particular, there are the “strides” and “padding”. Previously, we have used the simplest (but perhaps not most standard) setting of $\text{stride}=1$ and $\text{padding}=0$.

Strides give the size of the “jumps” as we slide the weights over the input image: see examples below. As such, $H_{\text{out}} \approx H_{\text{in}}/\text{stride}$ and $W_{\text{out}} \approx W_{\text{in}}/\text{stride}$.

Padding is a region around the input image that we fill with zeros: see dashed border around the inputs in the examples below. Without padding, the outputs will tend to be smaller than the inputs, even with $\text{stride}=1$. That happened above in Fig. 1.

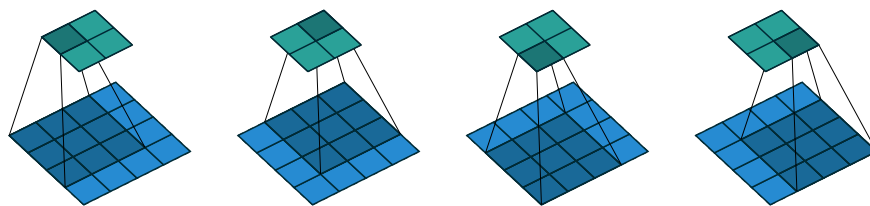


Figure 2: Convolving 3×3 weights over a 4×4 input using $\text{stride}=1$ and $\text{padding}=0$.

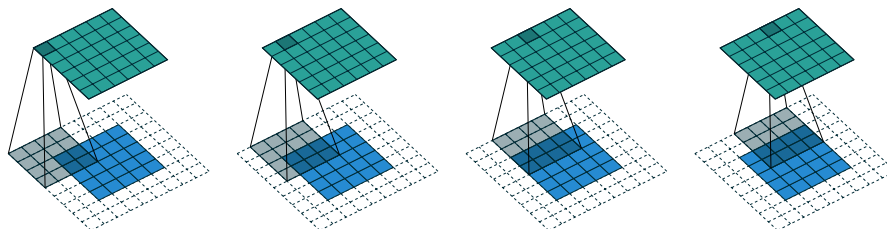


Figure 3: Convolving 4×4 weights over a 5×5 input with $\text{stride}=1$ and $\text{padding}=2$.

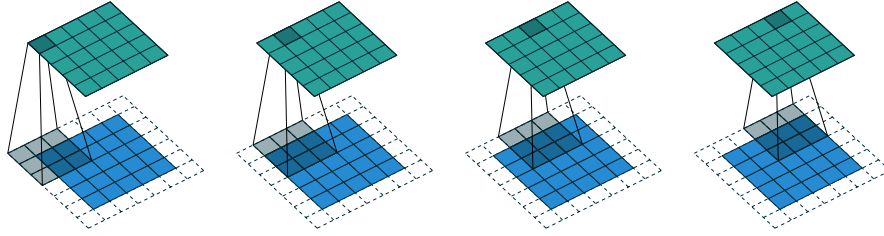


Figure 4: Convolution of 3×3 weights over a 5×5 input with stride=1 and padding=1. Note this is known as “same” padding, as the input is the same shape as the output.

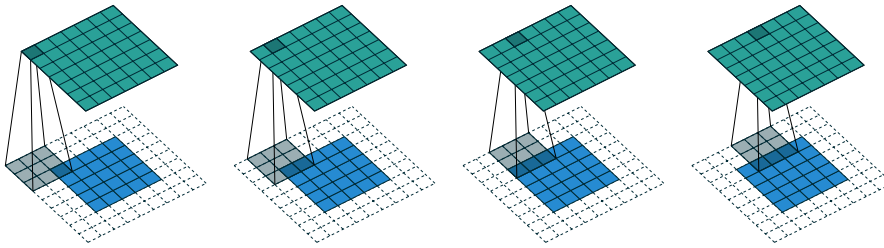


Figure 5: Convolution of 3×3 weights over a 5×5 input using stride=1 and padding=2.

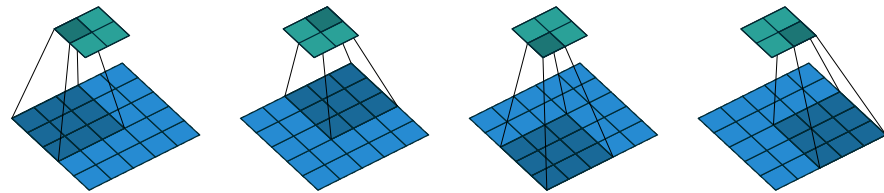


Figure 6: Convolution of 3×3 weights over a 5×5 input using stride=2 and padding=0.

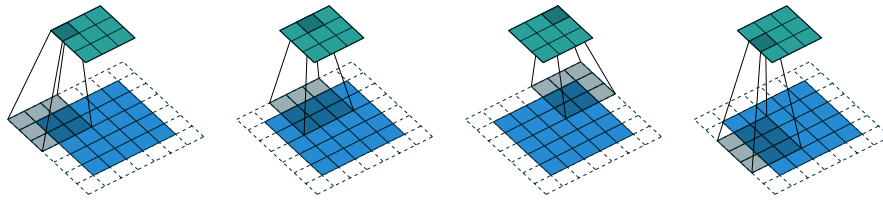


Figure 7: Convolving 3×3 weights over a 5×5 input with stride=2 and padding=1.

3 Pooling

In addition to discrete convolutions themselves, *pooling* operations are another important building block in CNNs. Pooling operations usually reduce the size of feature maps by using some function to summarize subregions, such as taking the average or the maximum value.

Pooling slides a window across the input, just like a convolution. The difference is that instead of multiplying the features in that window by weights, pooling layers feed the features in that layer into a simple *pooling function*. The pooling function is usually max or average. Again, pooling operations have strides and padding that behave in the same way as a standard convolution. Usually, when pooling we use a stride of at least 2, to ensure that we are actually reducing the size of the feature map.

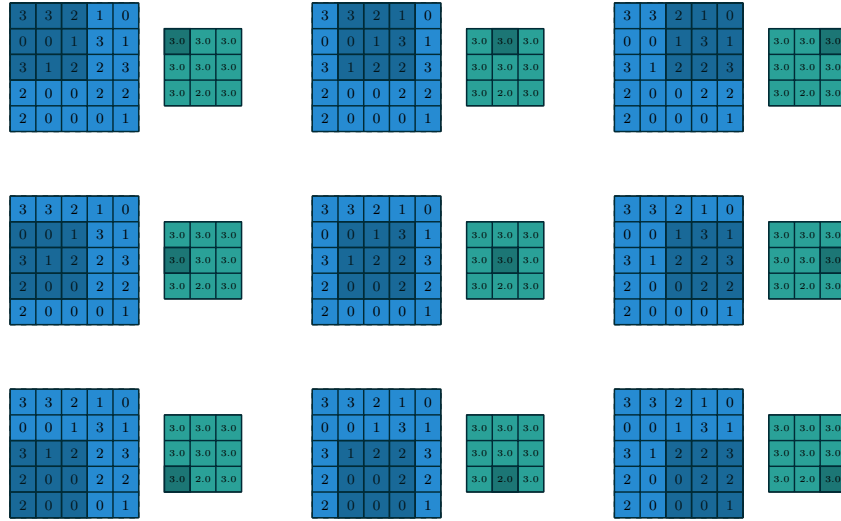


Figure 8: Computing the output values of a 3×3 max pooling operation on a 5×5 input using 1×1 strides.

3.1 Global average pooling

For image classification, we ultimately want just vector of logits as our output (for a single image). But convolutional layers are going to give us a whole feature map of size channels \times height \times width. So how can we convert a feature map into just a vector? Perhaps the most common method is “global average

pooling”, which just averages over the spatial dimensions.

$$\text{feature vector}_c = \frac{1}{HW} \sum_{x=1}^W \sum_{y=1}^H \text{feature map}_{c,x,y} \quad (17)$$

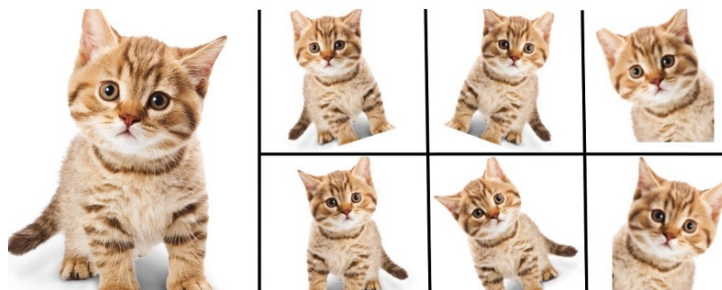
In PyTorch, it is easiest to implement this using `nn.AdaptiveAvgPool2d(1)`, which returns a feature map with height=1 and width=1.

4 Data Augmentation

Often I use convolutional neural networks for classification: i.e. for classifying what object is present in the image. Interestingly, it turns out that if I take an image of a kitten, and:

- Horizontal flip.
- Small rotation.
- Small translation/crop.

I still end up with an image of a kitten.



Enlarge your Dataset

So I can effectively enlarge my dataset by flipping/rotating/cropping all my images. And this turns out to be a *huge* factor in the success of modern networks. Basically all modern networks for images use some form of data augmentation.

5 Batchnorm

Batchnorm ensures that each channel has a specified, learned mean and variance. First, batchnorm computes the mean-and variance across channels, averaging across images, and spatial locations.

$$\mu_c = \frac{1}{NHW} \sum_{i,x,y} \text{input}_{i,c,x,y}, \quad (18)$$

$$\sigma_c^2 = \frac{1}{NHW} \sum_{i,x,y} (\text{input}_{i,c,x,y} - \mu_c)^2. \quad (19)$$

Then, batchnorm “standardises” the inputs: subtracting the per-channel mean, and dividing by the per-channel standard deviation,

$$\text{standardised input}_{i,c,x,y} = \frac{x_{i,c,x,y} - \mu_c}{\sigma_c}. \quad (20)$$

Importantly “standardised input” has mean zero and standard-deviation 1 for all channels. This is a bit restrictive: the neural network might perform better if the outputs have a different mean and standard deviation. Therefore, batchnorm finally multiplies by per-channel scale-factor, γ_c and adds a per-channel offset, β_c ,

$$\text{output}_{i,c,x,y} = \gamma_c \hat{x}_{i,c,x,y} + \beta_c. \quad (21)$$

Thus, the outputs of batchnorm have standard deviation γ_c and mean β_c .

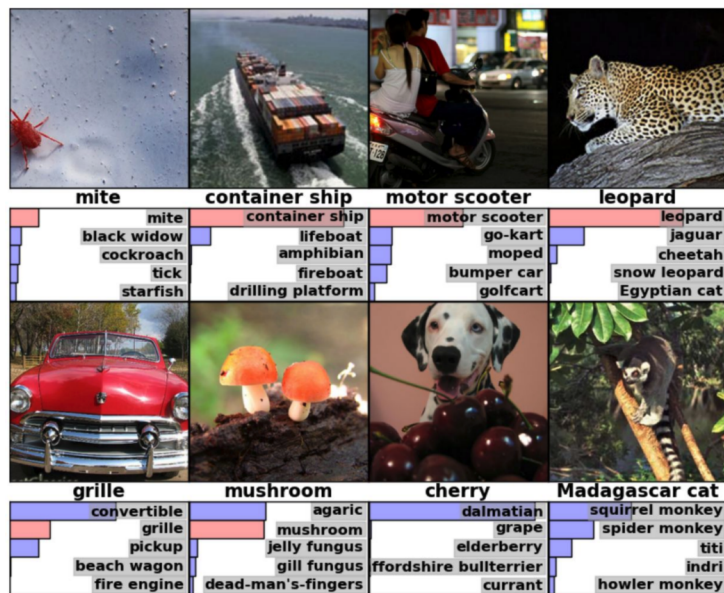
There are lots of hypotheses for how batchnorm works, but I haven’t really been convinced by any of them. Suffice it to say that batchnorm seems to stabilise training, and it seems to be *really* important for getting neural networks to work well. In fact, it gives perhaps the biggest boost to neural network performance of any tweak (except perhaps data augmentation, which also makes a huge difference).

6 Classic convolutional neural networks

We’re now in a position to understand the full architecture of some of the most classic convolutional neural networks (CNNs). Specifically, we’re going to look at arguably the two historically most important two neural networks: LeNet and AlexNet.

LeNet was developed in Yann LeCun et al. in 1998, and was designed to recognise images of handwritten digits (MNIST). This dataset included 60,000 greyscale, 28×28 images of handwritten digits. The resulting system was deployed in the real-world to read millions of checks (we have already seen this dataset). LeNet was the first demonstration that neural networks could do interesting tasks.

AlexNet was developed by Alex Krizhevsky et al. in 2012. If we have to choose a network that started the “deep learning revolution”, most people would choose AlexNet. AlexNet was designed to solve a much, much harder task: recognising one of 1000 objects in large (224×224 color images). Specifically, this dataset is known as ImageNet-1k (as it has 1000 object classes). This dataset has about one million training examples. AlexNet was about 50% better than the previous best method, based on hand-engineered features.



Lets zoom in a bit on the architecture:

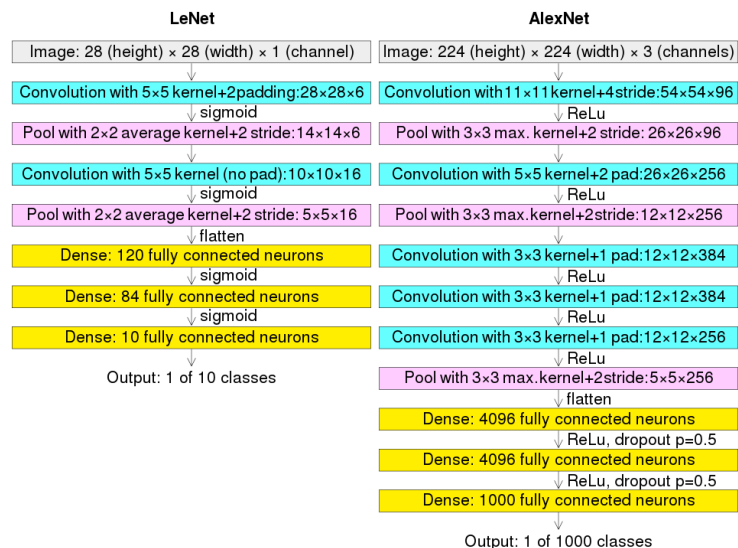


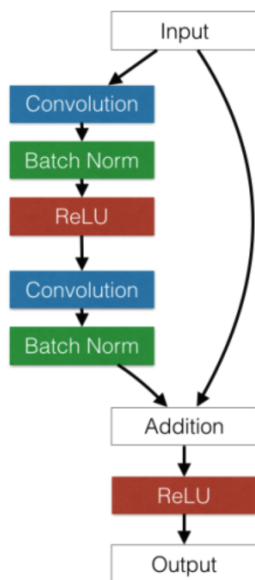
Figure 10: Each box is a fully-connected (dense), convolutional, or pooling layer. The size of the output is given in each box after the colon.

We can see the differences are surprisingly minimal. AlexNet's success is typically attributed to:

- Use of relu rather than sigmoid nonlinearities. As discussed previously, relu gives better backpropagation of gradients.
- AlexNet was one of the first times that GPUs were used to train neural networks, so they were able to use much more compute for training.
- Dropout for regularisation.
- An early variant of batchnorm after some of the nonlinearities (not shown on the diagram). This variant is no longer used, so is not discussed further in this course.

7 Modern convolutional neural networks: ResNets

Modern convolutional networks have a very similar structure. The only difference is the addition of a so-called “residual block”.



Several of these blocks would be strung together in a “ResNet”. Don’t worry about the precise details of exactly what order the blocks come in: there are *loads* of different variants for residual blocks anyway. The key idea is the “skip connection” on the right, which bypasses many layers. Again, this seems to work really well, but it isn’t really known why. One idea is that the skip connections help “signal” propagate through the network early on, and thereby help with training.